



Entwicklerhandbuch

AWS Data Pipeline



API-Version 2012-10-29

Copyright © 2025 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

AWS Data Pipeline: Entwicklerhandbuch

Copyright © 2025 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Die Handelsmarken und Handelsaufmachung von Amazon dürfen nicht in einer Weise in Verbindung mit nicht von Amazon stammenden Produkten oder Services verwendet werden, durch die Kunden irregeführt werden könnten oder Amazon in schlechtem Licht dargestellt oder diskreditiert werden könnte. Alle anderen Handelsmarken, die nicht Eigentum von Amazon sind, gehören den jeweiligen Besitzern, die möglicherweise zu Amazon gehören oder nicht, mit Amazon verbunden sind oder von Amazon gesponsert werden.

Table of Contents

.....	ix
Was ist AWS Data Pipeline?	1
Migration von Workloads von AWS Data Pipeline	2
Migration von Workloads zu AWS Glue	3
Migration von Workloads zu AWS Step Functions	4
Migration von Workloads zu Amazon MWAA	5
Kartierung der Konzepte	6
Beispiele	7
Zugehörige Services	8
Zugreifen AWS Data Pipeline	9
Preisgestaltung	10
Unterstützte Instance-Typen für Pipeline-Aktivitäten	10
EC2 Standard-Amazon-Instances nach AWS-Region	11
Zusätzliche unterstützte EC2 Amazon-Instances	12
Unterstützte EC2 Amazon-Instances für Amazon EMR-Cluster	13
AWS Data Pipeline Konzepte	15
Pipeline-Definition	15
Pipeline-Komponenten, Instances und Versuche	17
Task Runner	18
Datenknoten	19
Datenbanken	20
Aktivitäten	20
Vorbedingungen	21
Vom System verwaltete Vorbedingungen	22
Benutzerverwaltete Vorbedingungen	22
Ressourcen	22
Ressourcenlimits	23
Unterstützte Plattformen	23
Amazon EC2 Spot-Instances mit Amazon EMR-Clustern und AWS Data Pipeline	24
Aktionen	25
Proaktive Pipeline-Überwachung	26
Einrichten	27
Melden Sie sich an für AWS	27
Melden Sie sich für eine an AWS-Konto	27

Erstellen eines Benutzers mit Administratorzugriff	28
Erstellen Sie IAM-Rollen für AWS Data Pipeline und Pipeline-Ressourcen	29
Erlauben Sie IAM-Prinzipalen (Benutzern und Gruppen), die erforderlichen Aktionen auszuführen	30
Erteilen programmgesteuerten Zugriffs	31
Erste Schritte mit AWS Data Pipeline	34
Erstellen Sie die Pipeline	35
Überwachen der ausgeführten Pipeline	36
Anzeigen der Ausgabe	37
Löschen der Pipeline	37
Arbeiten mit Pipelines	38
Eine Pipeline erstellen	38
Erstellen Sie mit der CLI eine Pipeline aus Data Pipeline-Vorlagen	39
Anzeigen Ihrer Pipelines	59
Interpretieren der Pipeline-Statuscodes	59
Interpretieren des Pipeline- und Komponenten-Zustands	61
Anzeigen Ihrer Pipeline-Definitionen	63
Anzeigen von Pipeline-Instance Details	64
Anzeigen von Pipeline-Protokollen	64
Bearbeiten Ihrer Pipeline	66
Einschränkungen	66
Bearbeiten einer Pipeline mit dem AWS CLI	67
Klonen Ihrer Pipeline	68
Tagging Ihrer Pipeline	69
Deaktivieren Ihrer Pipeline	70
Deaktivieren Ihrer Pipeline über die AWS CLI	70
Löschen Ihrer Pipeline	71
Staging von Daten und Tabellen mit Aktivitäten	71
Daten-Staging mit ShellCommandActivity	73
Tabellen-Staging mit Hive und zum Staging fähigen Datenknoten	74
Tabellen-Staging mit Hive und nicht zum Staging fähigen Datenknoten	75
Verwenden von Ressourcen in mehreren Regionen	77
Cascading-Ausfälle und erneute Ausführungen	79
Aktivitäten	80
Datenknoten und Vorbedingungen	80
Ressourcen	80

Objekte, bei denen kaskadierte Fehler aufgetreten sind, erneut ausführen	80
Cascade-Fehler und Backfills	81
Syntax der Pipeline-Definitionsdatei	81
Dateistruktur	82
Pipeline-Felder	82
Benutzerdefinierte Felder	84
Arbeiten mit der API	85
Installieren des AWS-SDKs	85
Eine HTTP-Anfrage stellen an AWS Data Pipeline	86
Sicherheit	91
Datenschutz	92
Identitäts- und Zugriffsverwaltung	93
IAM-Richtlinien für AWS Data Pipeline	94
Beispielrichtlinien für AWS Data Pipeline	98
IAM-Rollen	101
Protokollieren und Überwachen	106
AWS Data Pipeline Informationen in CloudTrail	106
Grundlegendes zu AWS Data Pipeline Protokolldateieinträgen	107
Vorfallreaktion	109
Compliance-Validierung	109
Ausfallsicherheit	109
Infrastruktursicherheit	109
Konfiguration und Schwachstellenanalyse in AWS Data Pipeline	110
Tutorials	111
Daten mithilfe von Amazon EMR mit Hadoop Streaming verarbeiten	111
Bevor Sie beginnen	112
Verwenden der CLI	112
CSV-Daten von Amazon S3 zu Amazon S3 kopieren	117
Bevor Sie beginnen	118
Verwenden der CLI	119
Exportieren von MySQL-Daten nach Amazon S3	126
Bevor Sie beginnen	126
Verwenden der CLI	128
Daten nach Amazon Redshift kopieren	137
Bevor Sie beginnen: Konfigurieren Sie COPY-Optionen	138
Bevor Sie beginnen: Einrichten von Pipeline, Sicherheit und Cluster	139

Verwenden der CLI	140
Pipeline-Ausdrücke und -Funktionen	151
Einfache Datentypen	151
DateTime	151
Numerischer Wert	151
Objektverweise	151
Intervall	152
String	152
Ausdrücke	152
Verweisen auf Felder und Objekte	153
Verschachtelte Ausdrücke	154
Listen	155
Knotenausdruck	155
Ausdrucksauswertung	156
Mathematische Funktionen	157
Funktionen für Zeichenfolgen	157
Datums- und Zeitfunktionen	158
Sonderzeichen	167
Pipeline-Objektreferenz	169
Datenknoten	170
Dynamo-Knoten DBData	171
MySQLDataNode	178
RedshiftDataNode	186
S3 DataNode	194
SqlDataNode	202
Aktivitäten	209
CopyActivity	210
EmrActivity	218
HadoopActivity	228
HiveActivity	240
HiveCopyActivity	250
PigActivity	260
RedshiftCopyActivity	275
ShellCommandActivity	290
SqlActivity	301
Ressourcen	309

Ec2Resource	309
EmrCluster	321
HttpProxy	354
Vorbedingungen	357
Dynamo DBData ist vorhanden	357
Dynamo ist vorhanden DBTable	361
Vorhanden	365
S3 KeyExists	369
S3 PrefixNotEmpty	374
ShellCommandPrecondition	378
Datenbanken	383
JdbcDatabase	384
RdsDatabase	386
RedshiftDatabase	388
Datenformate	391
CSV-Datenformate	391
Custom Data Format	393
Dynamo-Format DBData	395
Dynamo DBExport DataFormat	398
RegEx Datenformat	400
TSV-Datenformate	402
Aktionen	404
SnsAlarm	404
Beenden	406
Plan	408
Beispiele	409
Syntax	413
Dienstprogramme	415
ShellScriptConfig	416
EmrConfiguration	417
Eigenschaft	422
Mit Task Runner arbeiten	426
Task Runner auf AWS Data Pipeline-Verwaltete Ressourcen	426
Arbeiten an vorhandenen Ressourcen mithilfe von Task Runner ausführen	428
Task Runner installieren	430
(Optional) Task Runner-Zugriff auf Amazon RDS gewähren	430

Task Runner starten	432
Die Task Runner-Protokollierung wird überprüft	433
Task Runner-Threads und Vorbedingungen	433
Task Runner-Konfigurationsoptionen	434
Task-Runner mit einem Proxy verwenden	437
Task Runner und Benutzerdefiniert AMIs	437
Fehlerbehebung	438
Suchen von Fehlern in Pipelines	438
Identifizieren des Amazon EMR-Clusters, der Ihre Pipeline bedient	439
Interpretieren der Pipeline-Statusdetails	440
Lokalisieren von Fehlerprotokollen	442
Pipeline-Protokolle	442
Hadoop-Job- und Amazon EMR-Schrittprotokolle	443
Beheben typischer Probleme	443
Pipeline bleibt im Status PENDING	443
Pipeline-Komponente bleibt im Status WAITING_FOR_RUNNER	444
Pipeline-Komponente bleibt im Status WAITING_ON_DEPENDENCIES	445
Ausführung beginnt nicht zum geplanten Zeitpunkt	446
Pipeline-Komponenten werden in der falschen Reihenfolge ausgeführt	446
EMR-Cluster schlägt mit Fehlermeldung fehl: The security token included in the request is invalid	447
Unzureichende Berechtigungen für den Zugriff auf Ressourcen	447
Statuscode: 400 Fehlercode: PipelineNotFoundException	447
Pipeline-Erstellung führt zu einem Sicherheits-Token-Fehler	447
Pipeline-Details werden nicht in der Konsole angezeigt	447
Error in remote runner Status Code: 404, AWS Service: Amazon S3	448
Access Denied - Not Authorized to Perform Function datapipeline:	448
Ältere Versionen von Amazon EMR AMIs können falsche Daten für große CSV-Dateien erstellen	449
AWS Data Pipeline Steigende Grenzwerte	449
Einschränkungen	450
Kontolimits	450
Limits für Webservice-Aufrufe	451
Überlegungen zur Skalierung	453
AWS Data Pipeline Ressourcen	454
Dokumentverlauf	456

AWS Data Pipeline ist für Neukunden nicht mehr verfügbar. Bestandskunden von AWS Data Pipeline können den Service weiterhin wie gewohnt nutzen. [Weitere Informationen](#)

Die vorliegende Übersetzung wurde maschinell erstellt. Im Falle eines Konflikts oder eines Widerspruchs zwischen dieser übersetzten Fassung und der englischen Fassung (einschließlich infolge von Verzögerungen bei der Übersetzung) ist die englische Fassung maßgeblich.

Was ist AWS Data Pipeline?

Note

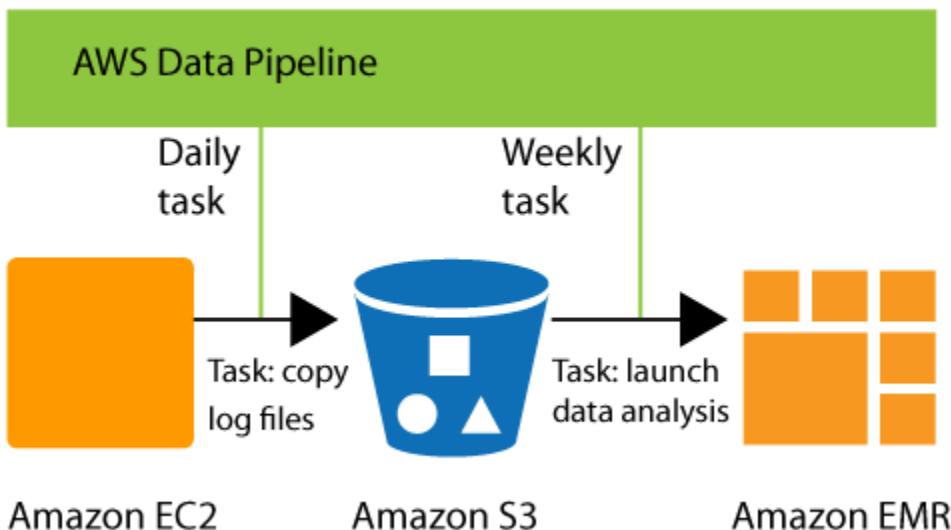
AWS Data Pipeline Der Service befindet sich im Wartungsmodus und es sind keine neuen Funktionen oder Regionserweiterungen geplant. Weitere Informationen und Informationen zur Migration Ihrer vorhandenen Workloads finden Sie unter [Migration von Workloads von AWS Data Pipeline](#)

AWS Data Pipeline ist ein Webservice, mit dem Sie die Bewegung und Transformation von Daten automatisieren können. Mit AWS Data Pipeline können Sie datengesteuerte Workflows definieren, sodass Aufgaben vom erfolgreichen Abschluss früherer Aufgaben abhängig sein können. Sie definieren die Parameter Ihrer Datentransformationen und setzen AWS Data Pipeline die von Ihnen eingerichtete Logik durch.

Die folgenden Komponenten von AWS Data Pipeline arbeiten zusammen, um Ihre Daten zu verwalten:

- Eine Pipeline-Definition legt die geschäftliche Logik der Datenverwaltung fest. Weitere Informationen finden Sie unter [Syntax der Pipeline-Definitionsdatei](#).
- Eine Pipeline plant und führt Aufgaben aus, indem sie EC2 Amazon-Instances zur Ausführung der definierten Arbeitsaktivitäten erstellt. Sie müssen nur die Pipeline-Definition in die Pipeline hochladen und diese anschließend aktivieren. Sie können auch die Pipeline-Definition einer gerade ausgeführten Pipeline bearbeiten. Sie müssen die Pipeline dann nur erneut aktivieren, damit die Änderungen wirksam werden. Außerdem können Sie die Pipeline deaktivieren, eine Datenquelle ändern und dann die Pipeline erneut aktivieren. Wenn Sie die Pipeline nicht mehr benötigen, können Sie sie löschen.
- Task Runner fragt nach Aufgaben ab und führt diese Aufgaben dann aus. Task Runner könnte zum Beispiel Protokolldateien in Amazon S3 kopieren und Amazon-EMR-Cluster starten. Task Runner ist installiert und wird automatisch auf Ressourcen ausgeführt, die durch Ihre Pipeline-Definitionen erstellt wurden. Sie können eine benutzerdefinierte Task Runner-Anwendung schreiben oder die Task Runner-Anwendung verwenden, die von bereitgestellt wird AWS Data Pipeline. Weitere Informationen finden Sie unter [Task Runner](#).

Sie können es beispielsweise verwenden, AWS Data Pipeline um die Protokolle Ihres Webservers jeden Tag in Amazon Simple Storage Service (Amazon S3) zu archivieren und dann einen wöchentlichen Amazon EMR-Cluster (Amazon EMR) über diese Protokolle laufen zu lassen, um Verkehrsberichte zu generieren. AWS Data Pipeline plant die täglichen Aufgaben zum Kopieren von Daten und die wöchentliche Aufgabe zum Starten des Amazon EMR-Clusters. AWS Data Pipeline stellt außerdem sicher, dass Amazon EMR wartet, bis die Daten des letzten Tages auf Amazon S3 hochgeladen wurden, bevor es mit der Analyse beginnt, auch wenn es zu unvorhergesehenen Verzögerungen beim Hochladen der Protokolle kommt.



Inhalt

- [Migration von Workloads von AWS Data Pipeline](#)
- [Zugehörige Services](#)
- [Zugreifen AWS Data Pipeline](#)
- [Preisgestaltung](#)
- [Unterstützte Instance-Typen für Pipeline-Aktivitäten](#)

Migration von Workloads von AWS Data Pipeline

AWS hat den AWS Data Pipeline Dienst 2012 eingeführt. Zu dieser Zeit waren Kunden auf der Suche nach einem Service, der ihnen hilft, Daten mithilfe einer Vielzahl von Rechenoptionen zuverlässig zwischen verschiedenen Datenquellen zu übertragen. Heute gibt es andere Dienste, die Kunden ein besseres Erlebnis bieten. Sie können zum Beispiel AWS Glue to verwenden, um Apache Spark-

Anwendungen auszuführen und zu orchestrieren, AWS Step Functions zur Orchestrierung von AWS Servicekomponenten oder Amazon Managed Workflows for Apache Airflow (Amazon MWAA), um die Workflow-Orchestrierung für Apache Airflow zu verwalten.

In diesem Thema wird erläutert, wie Sie von AWS Data Pipeline alternativen Optionen migrieren. Welche Option Sie wählen, hängt von Ihrem aktuellen Workload ab AWS Data Pipeline. Sie können typische Anwendungsfälle von entweder AWS Data Pipeline AWS Glue zu AWS Step Functions oder Amazon MWAA migrieren.

Migration von Workloads zu AWS Glue

[AWS Glue](#) ist ein Serverless-Datenintegrationsdienst, der es Analytics-Benutzern erleichtert, Daten aus mehreren Quellen zu erkennen, vorzubereiten, zu verschieben und zu integrieren. Es umfasst Tools für die Erstellung, Ausführung von Aufträgen und die Orchestrierung von Workflows. Mit AWS Glue können Sie mehr als 70 verschiedene Datenquellen entdecken und sich mit ihnen verbinden sowie Ihre Daten in einem zentralen Datenkatalog verwalten. Sie können ETL-Pipelines (Extract, Transform, Load) visuell erstellen, ausführen und überwachen, um Daten in Ihre Date Lakes zu laden. Außerdem können Sie mithilfe von Amazon Athena, Amazon EMR und Amazon Redshift Spectrum sofort katalogisierte Daten durchsuchen und abfragen.

Wir empfehlen, Ihren AWS Data Pipeline Workload zu den folgenden AWS Glue Zeiten zu migrieren:

- Sie suchen nach einem serverlosen Datenintegrationsservice, der verschiedene Datenquellen, Autorenschnittstellen wie visuelle Editoren und Notizbücher sowie erweiterte Datenverwaltungsfunktionen wie Datenqualität und Erkennung sensibler Daten unterstützt.
- Ihr Workload kann auf AWS Glue Workflows, Jobs (in Python oder Apache Spark) und Crawler migriert werden (Ihre bestehende Pipeline basiert beispielsweise auf Apache Spark).
- Sie benötigen eine einzige Plattform, die alle Aspekte Ihrer Datenpipeline abwickeln kann, einschließlich Aufnahme, Verarbeitung, Übertragung, Integritätstests und Qualitätsprüfungen.
- Ihre bestehende Pipeline wurde anhand einer vordefinierten Vorlage auf der AWS Data Pipeline Konsole erstellt, z. B. dem Exportieren einer DynamoDB-Tabelle nach Amazon S3, und Sie suchen nach derselben Vorlage für denselben Zweck.
- Ihre Arbeitslast hängt nicht von einer bestimmten Hadoop-Ökosystemanwendung wie Apache Hive ab.
- Ihr Workload erfordert keine Orchestrierung von lokalen Servern.

AWS berechnet einen sekundengenauen Stundensatz für Crawler (Datenermittlung) und ETL-Aufträge (Verarbeitung und Laden von Daten). AWS Glue Studio ist eine integrierte Orchestrierungs-Engine für AWS Glue Ressourcen, die ohne zusätzliche Kosten angeboten wird. Weitere Informationen zur Preisgestaltung finden Sie unter [AWS Glue Preise](#).

Migration von Workloads zu AWS Step Functions

[AWS Step Functions](#) ist ein Serverless-Orchestrierungsservice, mit dem Sie Workflows für Ihre geschäftskritischen Anwendungen erstellen können. Mit Step Functions verwenden Sie einen visuellen Editor, um Workflows zu erstellen und direkt in über 11.000 Aktionen für über 250 AWS Services wie AWS Lambda, Amazon EMR, DynamoDB und mehr zu integrieren. Sie können Step Functions verwenden, um Datenverarbeitungspipelines zu orchestrieren, Fehler zu behandeln und mit den Drosselungsgrenzen für die zugrunde liegenden Dienste zu arbeiten. Sie können Workflows erstellen, die Modelle für maschinelles Lernen verarbeiten und veröffentlichen, Mikroservices orchestrieren und AWS Dienste steuern, z. B. um ETL-Workflows (Extrahieren AWS Glue, Transformieren und Laden) zu erstellen. Sie können auch lang andauernde, automatisierte Workflows für Anwendungen erstellen, die menschliche Interaktion erfordern.

Ähnlich AWS Data Pipeline wie AWS Step Functions ist ein vollständig verwalteter Service von AWS. Sie müssen sich nicht um die Infrastruktur kümmern, Worker patchen, Betriebssystem-Versionsupdates oder Ähnliches verwalten.

Wir empfehlen, Ihren AWS Data Pipeline Workload zu AWS Step Functions zu migrieren, wenn:

- Sie suchen nach einem serverlosen, hochverfügbaren Workflow-Orchestrierungsservice.
- Sie suchen nach einer kostengünstigen Lösung, die mit der Granularität der Ausführung einer einzelnen Aufgabe berechnet wird.
- Ihre Workloads orchestrieren Aufgaben für mehrere andere AWS Services wie Amazon EMR, Lambda oder DynamoDB. AWS Glue
- Sie suchen nach einer Low-Code-Lösung, die über einen drag-and-drop visuellen Designer für die Workflow-Erstellung verfügt und für die Sie keine neuen Programmierkonzepte erlernen müssen.
- Sie suchen nach einem Service, der Integrationen mit über 250 anderen AWS Diensten bietet, die über 11.000 Aktionen abdecken out-of-the-box, sowie Integrationen mit benutzerdefinierten Nichtdiensten und Aktivitäten ermöglicht.

AWS Data Pipeline Sowohl Step Functions als auch Step Functions verwenden das JSON-Format, um Workflows zu definieren. Auf diese Weise können Sie Ihre Workflows in der Quellcodeverwaltung

speichern, Versionen verwalten, den Zugriff kontrollieren und mit CI/CD automatisieren. Step Functions verwendet eine Syntax namens Amazon State Language, die vollständig auf JSON basiert und einen nahtlosen Übergang zwischen der textuellen und der visuellen Darstellung des Workflows ermöglicht.

Mit Step Functions können Sie dieselbe Version von Amazon EMR auswählen, die Sie derzeit verwenden. AWS Data Pipeline

Für die Migration von Aktivitäten auf AWS Data Pipeline verwalteten Ressourcen können Sie die [AWS SDK-Serviceintegration](#) in Step Functions verwenden, um die Bereitstellung und Bereinigung von Ressourcen zu automatisieren.

Für die Migration von Aktivitäten auf lokalen Servern, benutzerverwalteten EC2 Instanzen oder einem benutzerverwalteten EMR-Cluster können Sie einen [SSM-Agent](#) auf der Instanz installieren. Sie können den Befehl über den [AWS Systems Manager Run Command](#) von Step Functions aus initiieren. Sie können den Zustandsmaschine auch über den in [Amazon](#) definierten Zeitplan initiieren EventBridge.

AWS Step Functions hat zwei Arten von Workflows: Standard-Workflows und Express-Workflows. Bei Standard-Workflows werden Ihnen Gebühren auf der Grundlage der Anzahl der Statusübergänge berechnet, die für die Ausführung Ihrer Anwendung erforderlich sind. Bei Express-Workflows werden Ihnen Gebühren auf der Grundlage der Anzahl der Anfragen für Ihren Workflow und dessen Dauer berechnet. Weitere Informationen zur Preisgestaltung finden Sie unter [AWS Step Functions Pricing](#).

Migration von Workloads zu Amazon MWAA

[Amazon MWAA](#) (Managed Workflows for Apache Airflow) ist ein verwalteter Orchestrierungs-Service für [Apache Airflow](#), der die Einrichtung und den Betrieb von end-to-end Datenpipelines in der Cloud im großen Maßstab erleichtert. Apache Airflow ist ein Open-Source-Tool, mit dem Abläufe und Aufgaben, die als „Workflows“ bezeichnet werden, programmgesteuert erstellt, geplant und überwacht werden können. Mit Amazon MWAA können Sie Airflow und die Programmiersprache Python verwenden, um Workflows zu erstellen, ohne die zugrunde liegende Infrastruktur im Hinblick auf Skalierbarkeit, Verfügbarkeit und Sicherheit verwalten zu müssen. Amazon MWAA passt seine Workflow-Ausführungskapazität automatisch an Ihre Bedürfnisse an und ist in AWS Sicherheitsservices integriert, um Ihnen einen schnellen und sicheren Zugriff auf Ihre Daten zu ermöglichen.

Ähnlich AWS Data Pipeline wie Amazon MWAA ist ein vollständig verwalteter Service, der von bereitgestellt wird. AWS Sie müssen sich zwar mit einigen neuen Konzepten vertraut machen, die für

diese Services spezifisch sind, aber Sie müssen sich nicht mit der Verwaltung der Infrastruktur, der Patch-Worker, der Verwaltung von Betriebssystemversions-Updates oder ähnlichem befassen.

Wir empfehlen, Ihre AWS Data Pipeline Workloads zu Amazon MWAA zu migrieren, wenn:

- Sie suchen nach einem verwalteten, hochverfügbaren Service zur Orchestrierung von in Python geschriebenen Workflows.
- Sie möchten auf eine vollständig verwaltete, weit verbreitete Open-Source-Technologie, Apache Airflow, umsteigen, um maximale Portabilität zu gewährleisten.
- Sie benötigen eine einzige Plattform, die alle Aspekte Ihrer Datenpipeline abwickeln kann, einschließlich Erfassung, Verarbeitung, Übertragung, Integritätstests und Qualitätsprüfungen.
- Sie suchen nach einem Service, der für die Orchestrierung von Daten-Pipelines konzipiert ist und Funktionen wie eine umfangreiche Benutzeroberfläche für Beobachtbarkeit, Neustarts für fehlgeschlagene Workflows, Backfills und Wiederholungsversuche für Aufgaben bietet.
- Sie sind auf der Suche nach einem Service, der über 800 vorgefertigte Bediener und Sensoren verfügt und AWS sowohl Dienstleistungen als auch andere Dienste abdeckt.

Amazon MWAA-Workflows sind mithilfe von Python als Directed Acyclic Graphs (DAGs) definiert, sodass Sie sie auch als Quellcode behandeln können. Mit dem erweiterbaren Python-Framework von Airflow können Sie Workflows erstellen, die sich mit praktisch jeder Technologie verbinden. Es verfügt über eine umfangreiche Benutzeroberfläche zum Anzeigen und Überwachen von Workflows und kann problemlos in Versionskontrollsysteme integriert werden, um den CI/CD-Prozess zu automatisieren.

Mit Amazon MWAA können Sie dieselbe Version von Amazon EMR wählen, die Sie derzeit verwenden. AWS Data Pipeline

AWS berechnet die Zeit, in der Ihre Airflow-Umgebung läuft, plus zusätzliche auto Skalierung, um mehr Mitarbeiter- oder Webserverkapazität bereitzustellen. Erfahren Sie mehr über die Preisgestaltung in [Amazon Managed Workflows for Apache Airflow Pricing](#).

Kartierung der Konzepte

Die folgende Tabelle enthält eine Übersicht der wichtigsten Konzepte, die von den Diensten verwendet werden. Es hilft Personen, die mit Data Pipeline vertraut sind, die Step Functions und die MWAA-Terminologie zu verstehen.

Data Pipeline	Glue	Step Functions	Amazon MWAA
Pipelines	Arbeitsabläufe	Arbeitsabläufe	Direkte Acrylgrafiken
Pipeline-Definition JSON	Workflow-Definition oder Python-basierte Blueprints	Amazon State Language JSON	Python-basiert
Aktivitäten	Aufträge	Staaten und Aufgaben	Aufgaben (Bediener und Sensoren)
Instances	Job läuft	Hinrichtungen	DAG läuft
Attempts	Versuche es erneut	Fänger und Retrier	Wiederholversuche
Pipeline-Zeitplan	Trigger einplanen	EventBridge Scheduler-Aufgaben	Cron, Zeitpläne, datenbewusst
Pipeline-Ausdrücke und -Funktionen	Blueprint-Bibliothek	Step Functions, intrinsische Funktionen und Lambda AWS	Erweiterbares Python-Framework

Beispiele

In den folgenden Abschnitten sind öffentliche Beispiele aufgeführt, auf die Sie zurückgreifen können, um von zu einzelnen Diensten AWS Data Pipeline zu migrieren. Sie können sie als Beispiele verwenden und Ihre eigene Pipeline für die einzelnen Dienste erstellen, indem Sie sie auf der Grundlage Ihres Anwendungsfalls aktualisieren und testen.

AWS Glue Beispiele

Die folgende Liste enthält Beispielimplementierungen für die häufigsten AWS Data Pipeline Anwendungsfälle mit AWS Glue

- [Spark-Jobs ausführen](#)
- [Daten von JDBC nach Amazon S3 kopieren \(einschließlich Amazon Redshift\)](#)
- [Daten aus Amazon S3 in JDBC kopieren \(einschließlich Amazon Redshift\)](#)
- [Daten von Amazon S3 nach DynamoDB kopieren](#)

- [Daten in und von Amazon Redshift verschieben](#)
- [Konten- und regionenübergreifender Zugriff auf DynamoDB-Tabellen](#)

AWS Beispiele für Step Functions

Die folgende Liste enthält Beispielimplementierungen für die häufigsten AWS Data Pipeline Anwendungsfälle mit AWS Step Functions.

- [Einen Amazon EMR-Job verwalten](#)
- [Ausführen eines Datenverarbeitungsauftrags auf Amazon EMR Serverless](#)
- [Jobs werden ausgeführt Hive/Pig/Hadoop](#)
- [Abfragen großer Datensätze \(Amazon Athena, Amazon S3,\) AWS Glue](#)
- [Ausführen von ETL-Workflows mit Amazon Redshift](#)
- [Crawler orchestrieren AWS Glue](#)

Sehen Sie sich zusätzliche [Tutorials](#) und [Beispielprojekte](#) zur Verwendung von AWS Step Functions an.

Amazon MWAA-Beispiele

Die folgende Liste enthält Beispielimplementierungen für die häufigsten AWS Data Pipeline Anwendungsfälle mit Amazon MWAA.

- [Einen Amazon EMR-Job ausführen](#)
- [Ein benutzerdefiniertes Plugin für Apache Hive und Hadoop erstellen](#)
- [Daten aus Amazon S3 in Redshift kopieren](#)
- [Ausführen eines Shell-Skripts auf einer Remote-Instance EC2](#)
- [Orchestrierung hybrider \(vor Ort\) Workflows](#)

Sehen Sie sich zusätzliche [Tutorials](#) und [Beispielprojekte](#) für die Verwendung von Amazon MWAA an.

Zugehörige Services

AWS Data Pipeline arbeitet mit den folgenden Services, um Daten zu speichern.

- Amazon DynamoDB — Bietet eine vollständig verwaltete NoSQL-Datenbank mit schneller Leistung zu niedrigen Kosten. Weitere Informationen finden Sie im [Amazon DynamoDB DynamoDB-Entwicklerhandbuch](#).
- Amazon RDS — Bietet eine vollständig verwaltete relationale Datenbank, die auf große Datenmengen skaliert werden kann. Weitere Informationen finden Sie im [Amazon Relational Database Service Developer Guide](#).
- Amazon Redshift — Bietet ein schneller, vollständig verwalteter Data Warehouse-Service für Datenmengen im Petabyte-Bereich, mit dem Sie große Datenmengen einfach und wirtschaftlich analysieren können. Weitere Informationen finden Sie im [Amazon Redshift Database Developer Guide](#).
- Amazon S3 — Bietet sicheren, dauerhaften und hoch skalierbaren Objektspeicher. Weitere Informationen finden Sie im [Amazon Simple Storage Service-Benutzerhandbuch](#).

AWS Data Pipeline arbeitet mit den folgenden Rechendiensten zusammen, um Daten zu transformieren.

- Amazon EC2 — bietet skalierbare Rechenkapazität — buchstäblich Server in den Rechenzentren von Amazon —, die Sie zum Aufbau und Hosten Ihrer Softwaresysteme verwenden. Weitere Informationen finden Sie im [EC2 Amazon-Benutzerhandbuch](#).
- Amazon EMR — Ermöglicht es Ihnen, große Datenmengen einfach, schnell und kostengünstig auf EC2 Amazon-Servern zu verteilen und zu verarbeiten, indem Sie ein Framework wie Apache Hadoop oder Apache Spark verwenden. Weitere Informationen finden Sie im [Amazon EMR Developer Guide](#).

Zugreifen AWS Data Pipeline

Sie können Ihre Pipelines über die folgenden Schnittstellen erstellen und verwalten:

- AWS-Managementkonsole— Bietet eine Webschnittstelle für den Zugriff AWS Data Pipeline.
- AWS Command Line Interface (AWS CLI) — Bietet Befehle für eine breite Palette von AWS-Services, einschließlich Windows, macOS und Linux AWS Data Pipeline, und wird unter Windows, macOS und Linux unterstützt. Weitere Informationen zur Installation von finden Sie unter [AWS Command Line Interface](#). AWS CLI Eine Liste der Befehle für AWS Data Pipeline finden Sie unter [datapipeline](#).

- AWS SDKs — Bietet sprachspezifische Funktionen APIs und übernimmt viele der Verbindungsdetails, wie zum Beispiel die Berechnung der Signaturen, die Verarbeitung des erneuten Absendens von Anforderungen und die Fehlerbehandlung. Weitere Informationen finden Sie unter [AWS SDKs](#).
- Abfrage-API — Bietet Low-Level-APIs APIs , die Sie mithilfe von HTTPS-Anforderungen aufrufen. Die Verwendung der Abfrage-API ist die direkteste Möglichkeit für den Zugriff auf AWS Data Pipeline. Allerdings müssen dann viele technische Abläufe, wie beispielsweise das Erzeugen des Hashwerts zum Signieren der Anforderung und die Fehlerbehandlung, in der Anwendung durchgeführt werden. Weitere Informationen finden Sie in der [AWS Data Pipeline -API-Referenz](#).

Preisgestaltung

Mit Amazon Web Services bezahlen Sie nur für das, was Sie tatsächlich nutzen. Denn AWS Data Pipeline Sie zahlen für Ihre Pipeline auf der Grundlage, wie oft und wo Ihre Aktivitäten und Voraussetzungen für die Ausführung geplant sind. Weitere Informationen finden Sie unter [AWS Data Pipeline – Preise](#).

Wenn Ihr AWS-Konto jünger als 12 Monate ist, sind Sie zur Nutzung des kostenlosen Kontingents berechtigt. Das kostenlose Kontingent umfasst drei Vorbedingungen mit geringer Häufigkeit und fünf Aktivitäten mit geringer Häufigkeit pro Monat. Weitere Informationen finden Sie unter [Kostenloses Kontingent für AWS](#).

Unterstützte Instance-Typen für Pipeline-Aktivitäten

Wenn eine Pipeline AWS Data Pipeline ausgeführt wird, kompiliert sie die Pipeline-Komponenten, um eine Reihe von umsetzbaren EC2 Amazon-Instances zu erstellen. Jede Instance enthält alle Informationen, die zum Ausführen einer bestimmten Aufgabe benötigt werden. Der komplette Satz an Instances stellt die To-do-Liste der Pipeline dar. AWS Data Pipeline übergibt die Instances zur Verarbeitung an Task Runner.

EC2 Instances haben verschiedene Konfigurationen, die als Instance-Typen bezeichnet werden. Jeder Instance-Typ verfügt über eine andere CPU, Eingabe/Ausgabe und Speicherkapazität. Zusätzlich zum Instance-Typ für eine Aktivität können Sie verschiedene Kaufoptionen auswählen. Nicht alle Instance-Typen stehen in allen AWS-Regionen zur Verfügung. Wenn ein Instance-Typ nicht verfügbar ist, kann Ihre Pipeline möglicherweise nicht bereitgestellt werden oder wird bei der Bereitstellung eingefroren. Informationen zur Instance-Verfügbarkeit finden Sie auf der [EC2 Amazon-Preisseite](#). Öffnen Sie den Link für Ihre Instance-Kaufoptionen und filtern Sie nach Region, um zu

sehen, ob ein Instance-Typ in dieser Region verfügbar ist. Weitere Informationen zu diesen Instance-Typen, Familien und Virtualisierungstypen finden Sie unter [Amazon EC2 Instances](#) und [Amazon Linux AMI Instance Type Matrix](#).

In den folgenden Tabellen werden die AWS Data Pipeline unterstützten Instance-Typen beschrieben. Sie können AWS Data Pipeline damit EC2 Amazon-Instances in jeder Region starten, auch in Regionen, in denen AWS Data Pipeline dies nicht unterstützt wird. Informationen zu Regionen, in denen dies unterstützt AWS Data Pipeline wird, finden Sie unter [AWS-Regionen und Endpunkte](#).

Inhalt

- [EC2 Standard-Amazon-Instances nach AWS-Region](#)
- [Zusätzliche unterstützte EC2 Amazon-Instances](#)
- [Unterstützte EC2 Amazon-Instances für Amazon EMR-Cluster](#)

EC2 Standard-Amazon-Instances nach AWS-Region

Wenn Sie in Ihrer Pipeline-Definition keinen Instance-Typ angeben, startet AWS Data Pipeline standardmäßig eine Instance.

In der folgenden Tabelle sind die EC2 Amazon-Instances aufgeführt, die standardmäßig in den Regionen AWS Data Pipeline verwendet werden, in denen dies unterstützt AWS Data Pipeline wird.

Name der Region	Region	Instance-Typ
USA Ost (Nord-Virginia)	us-east-1	m1.small
USA West (Oregon)	us-west-2	m1.small
Asien-Pazifik (Sydney)	ap-southeast-2	m1.small
Asien-Pazifik (Tokio)	ap-northeast-1	m1.small
EU (Irland)	eu-west-1	m1.small

In der folgenden Tabelle sind die EC2 Amazon-Instances aufgeführt, die standardmäßig in den Regionen AWS Data Pipeline gestartet werden, in denen AWS Data Pipeline dies nicht unterstützt wird.

Name der Region	Region	Instance-Typ
USA Ost (Ohio)	us-east-2	t2.small
USA West (Nordkalifornien)	us-west-1	m1.small
Asien-Pazifik (Mumbai)	ap-south-1	t2.small
Asien-Pazifik (Singapur)	ap-southeast-1	m1.small
Asien-Pazifik (Seoul)	ap-northeast-2	t2.small
Kanada (Zentral)	ca-central-1	t2.small
EU (Frankfurt)	eu-central-1	t2.small
EU (London)	eu-west-2	t2.small
EU (Paris)	eu-west-3	t2.small
Südamerika (São Paulo)	sa-east-1	m1.small

Zusätzliche unterstützte EC2 Amazon-Instances

Neben den Standard-Instances, die erstellt werden, wenn Sie in Ihrer Pipeline-Definition keinen Instance-Typ angeben, werden auch die folgenden Instances unterstützt.

In der folgenden Tabelle sind die EC2 Amazon-Instances aufgeführt, die AWS Data Pipeline unterstützt und erstellt werden können, sofern angegeben.

Instance-Klasse	Instance-Typen
Allgemeine Zwecke	t2.nano t2.micro t2.small t2.medium t2.large
Für Datenverarbeitung optimiert	c3.large c3.xlarge c3.2xlarge c3.4xlarge c3.8xlarge c4.large c4.xlarge c4.2xlarge c4.4xlarge c4.8xlarge c5.xlarge c5.9xlarge c5.2xlarge c5.4xlarge c5.9xlarge c5.18xlarge c5d.xlarge c5d.2xlarge c5d.4xlarge c5d.9xlarge c5d.18xlarge

Instance-Klasse	Instance-Typen
RAM-optimiert	m3.medium m3.large m3.xlarge m3.2xlarge m4.large m4.xlarge m4.2xlarge m4.4xlarge m4.10xlarge m4.16xlarge m5.xlarge m5.2xlarge m5.4xlarge m5.12xlarge m5.24xlarge m5d.xlarge m5d.2xlarge m5d.4xlarge m5d.12xlarge m5d.24xlarge r3.large r3.xlarge r3.2xlarge r3.4xlarge r3.8xlarge r4.large r4.xlarge r4.2xlarge r4.4xlarge r4.8xlarge r4.16xlarge
Speicheroptimiert	i2.xlarge i2.2xlarge i2.4xlarge i2.8xlarge hs1.8xlarge g2.2xlarge g2.8xlarge d2.xlarge d2.2xlarge d2.4xlarge d2.8xlarge

Unterstützte EC2 Amazon-Instances für Amazon EMR-Cluster

In dieser Tabelle sind die EC2 Amazon-Instances aufgeführt, die Amazon EMR-Cluster AWS Data Pipeline unterstützen und für diese erstellen können, sofern angegeben. Weitere Informationen finden Sie unter [Instance-Flotten konfigurieren](#) im Verwaltungshandbuch für Amazon EMR.

Instance-Klasse	Instance-Typen
Allgemeine Zwecke	m1.small m1.medium m1.large m1.xlarge m3.xlarge m3.2xlarge
Für Datenverarbeitung optimiert	c1.medium c1.xlarge c3.xlarge c3.2xlarge c3.4xlarge c3.8xlarge cc1.4xlarge cc2.8xlarge c4.large c4.xlarge c4.2xlarge c4.4xlarge c4.8xlarge c5.xlarge c5.9xlarge c5.2xlarge c5.4xlarge c5.9xlarge c5.18xlarge c5d.xlarge c5d.2xlarge c5d.4xlarge c5d.9xlarge c5d.18xlarge
RAM-optimiert	m2.xlarge m2.2xlarge m2.4xlarge r3.xlarge r3.2xlarge r3.4xlarge r3.8xlarge cr1.8xlarge m4.large m4.xlarge m4.2xlarge m4.4xlarge m4.10xlarge m4.16xlarge m5.xlarge m5.2xlarge m5.4xlarge m5.12xlarge m5.24xlarge m5d.xlarge m5d.2xlarge m5d.4xlarge m5d.12xlarge m5d.24xlarge

Instance-Klasse	Instance-Typen
	r4.large r4.xlarge r4.2xlarge r4.4xlarge r4.8xlarge r4.16xlarge
Speicheroptimiert	h1.4xlarge hs1.2xlarge hs1.4xlarge hs1.8xlarge i2.xlarge i2.2xlarge i2.4large i2.8xlarge d2.xlarge d2.2xlarge d2.4xlarge d2.8xlarge
Beschleunigtes Computing	g2.2xlarge cg1.4xlarge

AWS Data Pipeline Konzepte

Bevor Sie beginnen, informieren Sie sich über die wichtigsten Konzepte und Komponenten für AWS Data Pipeline.

Inhalt

- [Pipeline-Definition](#)
- [Pipeline-Komponenten, Instances und Versuche](#)
- [Task Runner](#)
- [Datenknoten](#)
- [Datenbanken](#)
- [Aktivitäten](#)
- [Vorbedingungen](#)
- [Ressourcen](#)
- [Aktionen](#)

Pipeline-Definition

Eine Pipeline-Definition ist die Art und Weise, wie Sie Ihre Geschäftslogik kommunizieren AWS Data Pipeline. Sie umfasst die folgenden Informationen:

- Namen, Speicherorte und Formate der Datenquellen
- Aktivitäten, mit denen die Daten transformiert werden
- Den Zeitplan für diese Aktivitäten
- Ressourcen für die Ausführung Ihrer Aktivitäten und Vorbedingungen
- Voraussetzungen, die erfüllt werden müssen, bevor die Aktivitäten geplant werden können
- Möglichkeiten zur Information über Statusänderungen beim Fortschreiten der Pipeline-Durchführung

AWS Data Pipeline Ermittelt anhand Ihrer Pipeline-Definition die Aufgaben, plant sie und weist sie Task-Runnern zu. Wenn eine Aufgabe nicht erfolgreich abgeschlossen wird, AWS Data Pipeline versucht die Aufgabe gemäß Ihren Anweisungen erneut und weist sie gegebenenfalls einem anderen

Task-Runner zu. Sie können die Pipeline so konfigurieren, dass Sie benachrichtigt werden, wenn die Aufgabe wiederholt fehlschlägt.

In Ihrer Pipeline-Definition können Sie beispielsweise angeben, dass von Ihrer Anwendung generierte Protokolldateien 2013 jeden Monat in einem Amazon S3 S3-Bucket archiviert werden. AWS Data Pipeline würde dann 12 Aufgaben erstellen, von denen jede Daten eines Monats kopiert, unabhängig davon, ob der Monat 30, 31, 28 oder 29 Tage umfasste.

Eine Pipeline-Definition können Sie folgendermaßen erstellen:

- Grafisch, mithilfe der Konsole AWS Data Pipeline
- Textlich, indem Sie eine JSON-Datei in dem von der Befehlszeile verwendeten Format erstellen
- Programmgesteuert, indem der Webservice entweder mit einer der AWS SDKs - oder der API aufgerufen wird AWS Data Pipeline

Eine Pipeline-Definition kann die folgenden Komponententypen enthalten.

Pipeline-Komponenten

Datenknoten

Den Speicherort von Eingabedaten für eine Aufgabe oder der Speicherort, an dem die Ausgabedaten gespeichert werden sollen.

Aktivitäten

Eine Definition der Arbeit, die nach einem Zeitplan mit einer Datenverarbeitungsressource und (typischerweise) Eingabe- und Ausgabeknoten durchgeführt werden soll.

Vorbedingungen

Eine Bedingungsaussage, die erfüllt sein muss, damit eine Aktion ausgeführt werden kann.

Ressourcen

Die Datenverarbeitungsressource, die die Arbeit ausführt, die eine Pipeline definiert.

Aktionen

Eine Aktion, die ausgelöst wird, wenn bestimmte Bedingungen erfüllt sind, z. B. wenn eine Aktivität fehlschlägt.

Weitere Informationen finden Sie unter [Syntax der Pipeline-Definitionsdatei](#).

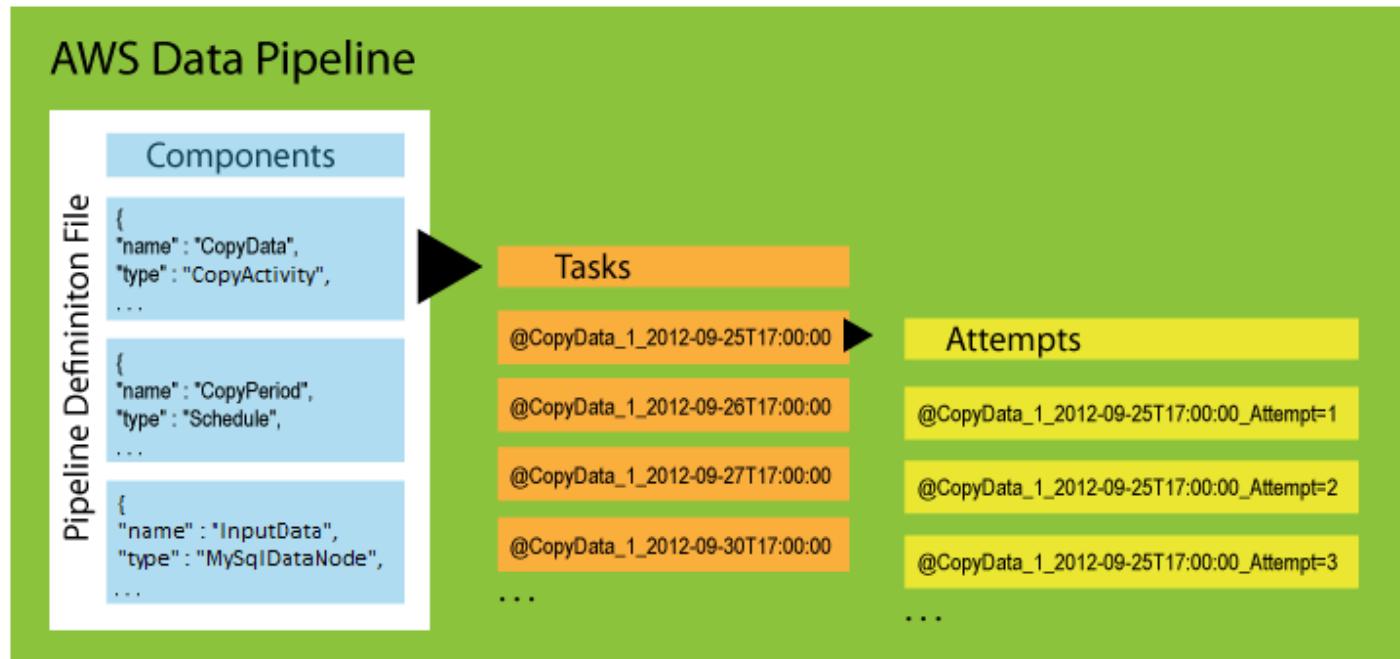
Pipeline-Komponenten, Instances und Versuche

Es gibt drei Komponententypen im Zusammenhang mit einer geplanten Pipeline:

- Pipeline-Komponenten — Pipeline-Komponenten stellen die Geschäftslogik der Pipeline dar und werden durch die verschiedenen Abschnitte einer Pipeline-Definition repräsentiert. Pipeline-Komponenten geben die Datenquellen, Aktivitäten, den Zeitplan und die Vorbedingungen des Workflows an. Sie können Eigenschaften von übergeordneten Komponenten übernehmen. Beziehungen zwischen Komponenten werden durch Verweise definiert. Pipeline-Komponenten definieren die Regeln für die Datenverwaltung.
- Instanzen — Wenn eine Pipeline AWS Data Pipeline ausgeführt wird, kompiliert sie die Pipeline-Komponenten, um eine Reihe von Instanzen zu erstellen, für die Aktionen erforderlich sind. Jede Instance enthält alle Informationen, die zum Ausführen einer bestimmten Aufgabe benötigt werden. Der vollständige Satz von Instanzen ist die Aufgabenliste der Pipeline. AWS Data Pipeline übergibt die Instanzen zur Bearbeitung an Task-Runner.
- Versuche – Um eine robuste Datenverwaltung sicherzustellen, wiederholt AWS Data Pipeline fehlgeschlagene Vorgänge. Diese Wiederholungen werden durchgeführt, bis die maximal erlaubte Anzahl an Wiederholungsversuchen erreicht ist. Versuchsobjekte verfolgen die einzelnen Versuche, Ergebnisse und ggf. Fehlergründe nach. Im Grunde ist es die Instanz mit einem Zähler. AWS Data Pipeline führt Wiederholungsversuche mit denselben Ressourcen wie bei den vorherigen Versuchen durch, z. B. Amazon EMR-Cluster und EC2 -Instances.

Note

Das Wiederholen fehlgeschlagener Aufgaben ist ein wichtiger Bestandteil einer Fehlertoleranzstrategie. AWS Data Pipeline -Definitionen stellen Bedingungen und Schwellenwerte zur Steuerung der erneuten Versuche bereit. Zu viele erneute Versuche können jedoch dazu führen, dass unwiederbringliche Fehler zu spät erkannt werden, da AWS Data Pipeline Fehler erst dann meldet, wenn die festgelegte Anzahl an Wiederholungsversuchen erreicht wurde. Diese zusätzlichen Wiederholungen können für zusätzliche Gebühren sorgen, wenn sie auf AWS-Ressourcen ausgeführt werden. Überlegen Sie sich daher sorgfältig, wann es angemessen ist, die AWS Data Pipeline Standardeinstellungen, die Sie zur Steuerung von Wiederholungsversuchen und ähnlichen Einstellungen verwenden, zu überschreiten.

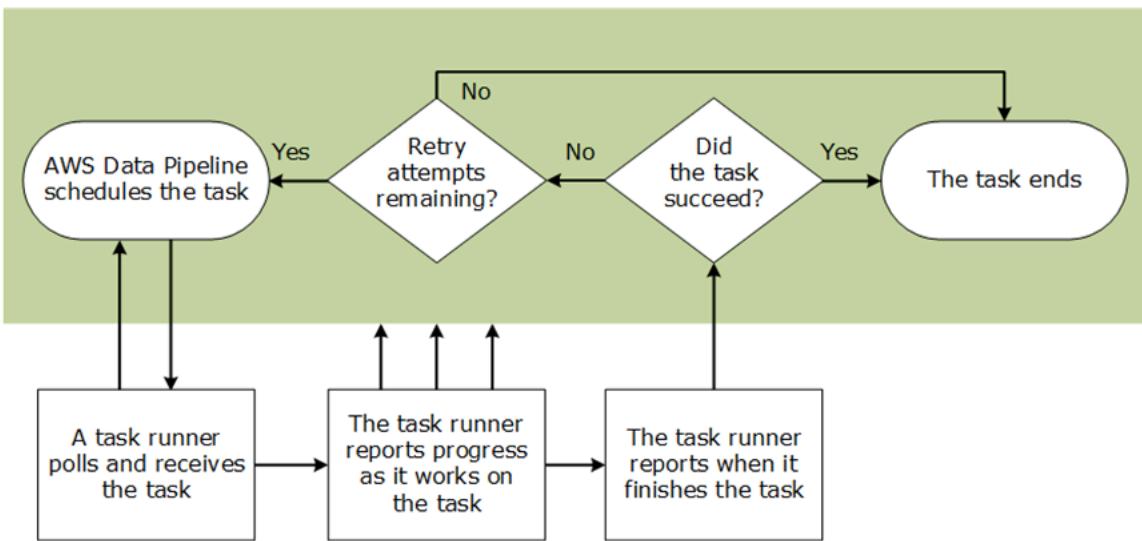


Task Runner

Ein Task-Runner ist eine Anwendung, die AWS Data Pipeline nach Aufgaben fragt und diese Aufgaben dann ausführt.

Task Runner ist eine Standardimplementierung eines Task Runners, der von AWS Data Pipeline bereitgestellt wird. Wenn Task Runner installiert und konfiguriert ist, fragt er AWS Data Pipeline nach Aufgaben ab, die mit Pipelines verknüpft sind, die Sie aktiviert haben. Wenn Task Runner eine Aufgabe zugewiesen wird, führt er diese Aufgabe aus und meldet ihren Status an. AWS Data Pipeline

Das folgende Diagramm zeigt, wie AWS Data Pipeline und wie ein Task-Runner interagieren, um eine geplante Aufgabe zu verarbeiten. Eine Aufgabe ist eine separate Arbeitseinheit, die sich der AWS Data Pipeline Service mit einem Task-Runner teilt. Sie unterscheidet sich von einer Pipeline, die eine allgemeine Definition von Aktivitäten und Ressourcen ist, die normalerweise zu mehreren Aufgaben führt.



Es gibt zwei Möglichkeiten, wie Sie Task Runner verwenden können, um Ihre Pipeline zu verarbeiten:

- AWS Data Pipeline installiert Task Runner für Sie auf Ressourcen, die vom AWS Data Pipeline Webservice gestartet und verwaltet werden.
- Sie installieren Task Runner auf einer Rechenressource, die Sie verwalten, z. B. auf einer EC2 Instanz mit langer Laufzeit oder auf einem lokalen Server.

Weitere Informationen zur Arbeit mit Task Runner finden Sie unter [Arbeiten mit Task Runner](#).

Datenknoten

In AWS Data Pipeline definiert ein Datenknoten den Speicherort und die Art der Daten, die eine Pipeline-Aktivität als Eingabe oder Ausgabe verwendet. AWS Data Pipeline unterstützt die folgenden Arten von Datenknoten:

[Dynamo-Knoten DBData](#)

Eine DynamoDB-Tabelle, die Daten für [HiveActivity](#) oder [EmrActivity](#) zur Verwendung enthält.

[SqlDataNode](#)

Eine SQL-Tabellen- und Datenbankabfrage, die Daten zur Verwendung durch eine Pipeline-Aktivität repräsentiert.

Note

Wurde zuvor verwendet MySqlDataNode . Verwenden Sie SqlDataNode stattdessen.

RedshiftDataNode

Eine Amazon Redshift Redshift-Tabelle, die Daten [RedshiftCopyActivity](#) zur Verwendung enthält.

S3 DataNode

Ein Amazon S3 S3-Speicherort, der eine oder mehrere Dateien enthält, die für eine Pipeline-Aktivität verwendet werden können.

Datenbanken

AWS Data Pipeline unterstützt die folgenden Datenbanktypen:

JdbcDatabase

Eine JDBC-Datenbank.

RdsDatabase

Eine Amazon RDS-Datenbank.

RedshiftDatabase

Eine Amazon Redshift Redshift-Datenbank.

Aktivitäten

In AWS Data Pipeline ist eine Aktivität eine Pipeline-Komponente, die die auszuführende Arbeit definiert. AWS Data Pipeline bietet mehrere vorkonfigurierte Aktivitäten, die allgemeine Szenarien berücksichtigen, z. B. das Verschieben von Daten von einem Speicherort an einen anderen, das Ausführen von Hive-Abfragen usw. Aktivitäten sind erweiterbar, sodass Sie Ihre eigenen benutzerdefinierten Skripts ausführen und so unzählige Kombinationen unterstützen können.

AWS Data Pipeline unterstützt die folgenden Arten von Aktivitäten:

[CopyActivity](#)

Kopiert Daten von einem Speicherort zu einem anderen.

[EmrActivity](#)

Führt einen Amazon EMR-Cluster aus.

[HiveActivity](#)

Führt eine Hive-Abfrage auf einem Amazon EMR-Cluster aus.

[HiveCopyActivity](#)

Führt eine Hive-Abfrage auf einem Amazon EMR-Cluster mit Unterstützung für erweiterte Datenfilterung und Unterstützung für [S3 DataNode](#) und aus. [Dynamo-Knoten DBData](#)

[PigActivity](#)

Führt ein Pig-Skript auf einem Amazon EMR-Cluster aus.

[RedshiftCopyActivity](#)

Kopiert Daten in und aus Amazon Redshift Redshift-Tabellen.

[ShellCommandActivity](#)

Führt einen benutzerdefinierten UNIX/Linux-Shell-Befehl als Aktivität aus.

[SqlActivity](#)

Führt eine SQL-Abfrage auf einer Datenbank aus.

Einige Aktivitäten bieten spezielle Unterstützung für Staging-Daten und Datenbanktabellen. Weitere Informationen finden Sie unter [Staging von Daten und Tabellen mit Pipeline-Aktivitäten](#).

Vorbedingungen

In ist eine Vorbedingung eine Pipeline-Komponente AWS Data Pipeline, die bedingte Anweisungen enthält, die erfüllt sein müssen, bevor eine Aktivität ausgeführt werden kann. Eine Vorbedingung kann beispielsweise überprüfen, ob Quelldaten vorhanden sind, bevor eine Pipelineaktivität versucht, sie zu kopieren. AWS Data Pipeline bietet mehrere vorkonfigurierte Vorbedingungen, die allgemeine Szenarien berücksichtigen, z. B. ob eine Datenbanktabelle existiert, ob ein Amazon S3

S3-Schlüssel vorhanden ist usw. Vorbedingungen sind jedoch erweiterbar, sodass Sie Ihre eigenen benutzerdefinierten Skripts ausführen und so unzählige Kombinationen unterstützen können.

Es gibt zwei Arten von Vorbedingungen: vom System verwaltete Vorbedingungen und benutzerverwaltete Vorbedingungen. Systemverwaltete Vorbedingungen werden vom AWS Data Pipeline Webservice in Ihrem Namen ausgeführt und erfordern keine Rechenressourcen. Benutzerverwaltete Vorbedingungen werden nur auf der Datenverarbeitungsressource ausgeführt, die Sie im Feld `runsOn` oder `workerGroup` festgelegt haben. Die `workerGroup`-Ressource ist von der Aktivität abgeleitet, die die Vorbedingung nutzt.

Vom System verwaltete Vorbedingungen

[Dynamo DBData ist vorhanden](#)

Prüft, ob Daten in einer bestimmten DynamoDB-Tabelle vorhanden sind.

[Dynamo ist vorhanden DBTable](#)

Prüft, ob eine DynamoDB-Tabelle existiert.

[S3 KeyExists](#)

Prüft, ob ein Amazon S3 S3-Schlüssel existiert.

[S3 PrefixNotEmpty](#)

Prüft, ob ein Amazon S3 S3-Präfix leer ist.

Benutzerverwaltete Vorbedingungen

[Vorhanden](#)

Prüft, ob ein Datenknoten vorhanden ist.

[ShellCommandPrecondition](#)

Führt einen Unix-/Linux-Shell-Befehl als Vorbedingung aus.

Ressourcen

Bei AWS Data Pipeline einer Ressource handelt es sich um die Rechenressource, die die in einer Pipelineaktivität festgelegte Arbeit ausführt. AWS Data Pipeline unterstützt die folgenden Ressourcentypen:

[Ec2Resource](#)

Eine EC2 Instanz, die die durch eine Pipeline-Aktivität definierte Arbeit ausführt.

[EmrCluster](#)

Ein Amazon EMR-Cluster, der die durch eine Pipeline-Aktivität definierte Arbeit ausführt, z. B.

[EmrActivity](#)

Ressourcen können in derselben Region mit ihrem Arbeitsdatensatz ausgeführt werden, auch in einer anderen Region als AWS Data Pipeline. Weitere Informationen finden Sie unter [Verwenden einer Pipeline mit Ressourcen in mehreren Regionen](#).

Ressourcenlimits

AWS Data Pipeline lässt sich skalieren, um eine große Anzahl gleichzeitiger Aufgaben zu bewältigen, und Sie können ihn so konfigurieren, dass automatisch die Ressourcen bereitgestellt werden, die für die Bewältigung großer Workloads erforderlich sind. Diese automatisch erstellten Ressourcen sind von Ihnen steuerbar und zählen bei Ihren Ressourcenlimits für Ihr AWS-Konto mit. Wenn Sie beispielsweise so konfigurieren, AWS Data Pipeline dass ein Amazon EMR-Cluster mit 20 Knoten automatisch zur Verarbeitung von Daten erstellt wird und Ihr AWS-Konto ein EC2 Instance-Limit von 20 hat, können Sie versehentlich Ihre verfügbaren Backfill-Ressourcen erschöpfen. Daher sollten Sie diese Ressourceneinschränkungen bei Ihrem Design berücksichtigen oder Ihre Kontolimits entsprechend erweitern. Weitere Informationen zu Service Limits finden Sie unter [AWS Service Limits](#) in der Allgemeinen AWS-Referenz.

Note

Der Grenzwert ist eine Instance pro Ec2Resource Komponentenobjekt.

Unterstützte Plattformen

Pipelines können Ihre Ressourcen in den folgenden Plattformen starten:

EC2-Klassiker

Ihre Ressourcen werden in einem einzelnen, flachen Netzwerk ausgeführt, das Sie gemeinsam mit anderen Kunden verwenden.

EC2-VPC

Ihre Ressourcen werden in einer Virtual Private Cloud (VPC) ausgeführt, die logisch von Ihrem AWS-Konto isoliert ist.

Ihr AWS-Konto kann Ressourcen auf regionaler Basis entweder auf beiden Plattformen oder nur in EC2 VPC. Weitere Informationen finden Sie unter [Unterstützte Plattformen](#) im EC2 Amazon-Benutzerhandbuch.

Wenn Ihr AWS-Konto nur EC2 -VPC unterstützt, erstellen wir in jeder AWS-Region eine Standard-VPC für Sie. Standardmäßig starten wir Ihre Ressourcen in einem Standard-Subnetz Ihrer Standard-VPC. Alternativ können Sie eine nicht standardmäßige VPC erstellen und bei der Konfiguration Ihrer Ressourcen eines ihrer Subnetze angeben. Dann starten wir Ihre Ressourcen in dem angegebenen Subnetz der nicht standardmäßigen VPC.

Wenn Sie eine Instance in einer VPC starten, müssen Sie eine Sicherheitsgruppe angeben, die speziell für diese VPC erstellt wurde. Sie können keine Sicherheitsgruppe angeben, die Sie für EC2 -Classic erstellt haben, wenn Sie eine Instance in einer VPC starten. Darüber hinaus müssen Sie die Sicherheitsgruppen-ID und nicht den Sicherheitsgruppennamen verwenden, um eine Sicherheitsgruppe für eine VPC festzulegen.

Amazon EC2 Spot-Instances mit Amazon EMR-Clustern und AWS Data Pipeline

Pipelines können Amazon EC2 Spot-Instances für die Task-Knoten in ihren Amazon EMR-Cluster-Ressourcen verwenden. Pipelines verwenden standardmäßig On-Demand-Instances. Mit Spot-Instances können Sie EC2 Ersatz-Instances verwenden und ausführen. Das Spot Instance-Preismodell ergänzt das On-Demand-Preismodell und das Reserved Instance-Preismodell und stellt womöglich die kosteneffizienteste Option für Rechenkapazität dar, je nach Anwendung. Weitere Informationen finden Sie auf der Produktseite von [Amazon EC2 Spot Instances](#).

Wenn Sie Spot-Instances verwenden, AWS Data Pipeline übermittelt Amazon EMR Ihren Spot-Instance-Höchstpreis, wenn Ihr Cluster gestartet wird. Es ordnet die Arbeit des Clusters automatisch der Anzahl der Spot-Instance-Task-Knoten zu, die Sie mithilfe des Felds definieren. `taskInstanceCount` AWS Data Pipeline schränkt Spot-Instances für Task-Nodes ein, um sicherzustellen, dass On-Demand-Core-Knoten für den Betrieb Ihrer Pipeline verfügbar sind.

Sie können eine fehlgeschlagene oder abgeschlossene Pipeline-Ressourcen-Instance bearbeiten und Spot-Instances hinzufügen. Wenn die Pipeline den Cluster erneut startet, nutzt er die Spot-Instances für die Aufgabenknoten.

Überlegungen zu Spot-Instances

Wenn Sie Spot-Instances mit verwenden AWS Data Pipeline, gelten die folgenden Überlegungen:

- Ihre Spot-Instances können beendet werden, wenn der Spot-Instance-Preis Ihren Höchstpreis für die Instance übersteigt oder wenn Amazon EC2 Kapazitätsgründe hat. Sie verlieren Ihre Daten jedoch nicht, da Cluster mit Kernknoten AWS Data Pipeline verwendet werden, bei denen es sich immer um On-Demand-Instances handelt, die nicht gekündigt werden können.
- Spot-Instances können mehr Zeit zum Starten benötigen, weil sie Kapazität asynchron bereitstellen. Aus diesem Grund läuft eine Spot-Instance-Pipeline möglicherweise langsamer als eine äquivalente On-Demand-Instance-Pipeline.
- Ihr Cluster wird möglicherweise nicht ausgeführt, wenn Sie Ihre Spot-Instances nicht erhalten, beispielsweise, wenn Ihr Höchstpreis zu niedrig ist.

Aktionen

AWS Data Pipeline Aktionen sind Schritte, die eine Pipeline-Komponente unternimmt, wenn bestimmte Ereignisse eintreten, wie z. B. erfolgreiche, fehlgeschlagene oder verspätete Aktivitäten. Das Ereignisfeld einer Aktivität bezieht sich auf eine Aktion, z. B. einen Verweis auf `snsalarm` im Feld `onLateAction` von `EmrActivity`.

AWS Data Pipeline verwendet Amazon SNS SNS-Benachrichtigungen als primäres Mittel, um den Status von Pipelines und ihren Komponenten unbeaufsichtigt anzuzeigen. Weitere Informationen finden Sie unter [Amazon SNS](#). Zusätzlich zu SNS-Benachrichtigungen können Sie die AWS Data Pipeline Konsole und die CLI verwenden, um Informationen zum Pipeline-Status abzurufen.

AWS Data Pipeline unterstützt die folgenden Aktionen:

SnsAlarm

Eine Aktion, die eine SNS-Benachrichtigung an ein Thema sendet, basierend auf den Ereignissen `onSuccess`, `OnFail` und `onLateAction`.

Beenden

Eine Aktion, die eine Stornierung von ausstehenden oder nicht abgeschlossenen Aktivitäten, Ressourcen oder Datenknoten auslöst. Sie können keine Aktionen beenden, die `onSuccess`, `OnFail` oder `onLateAction` beinhalten.

Proaktive Pipeline-Überwachung

Die beste Möglichkeit zum Erkennen von Problemen ist die proaktive Überwachung Ihrer Pipelines von Anfang an. Sie können Pipeline-Komponenten so konfigurieren, dass sie Sie über bestimmte Situationen oder Ereignisse informieren, z. B. wenn eine Pipeline-Komponente ausfällt oder nicht zur geplanten Startzeit startet. AWS Data Pipeline macht es einfach, Benachrichtigungen zu konfigurieren, indem es Ereignisfelder für Pipeline-Komponenten bereitstellt, die Sie mit Amazon SNS SNS-Benachrichtigungen verknüpfen können, wie `onSuccess``onFail`, `undonLateAction`.

Einrichtung für AWS Data Pipeline

Führen Sie vor der ersten Verwendung AWS Data Pipeline die folgenden Aufgaben aus.

Aufgaben

- [Melden Sie sich an für AWS](#)
- [Erstellen Sie IAM-Rollen für AWS Data Pipeline und Pipeline-Ressourcen](#)
- [Erlauben Sie IAM-Prinzipalen \(Benutzern und Gruppen\), die erforderlichen Aktionen auszuführen](#)
- [Erteilen programmgesteuerten Zugriffs](#)

Nachdem Sie diese Aufgaben abgeschlossen haben, können Sie mit der Verwendung beginnen AWS Data Pipeline. Ein Tutorial zum Einstieg finden Sie unter [Erste Schritte mit AWS Data Pipeline](#).

Melden Sie sich an für AWS

Wenn Sie sich für Amazon Web Services (AWS) registrieren, wird Ihr AWS-Konto automatisch für alle Services in AWS registriert, einschließlich AWS Data Pipeline. Berechnet werden Ihnen aber nur die Services, die Sie nutzen. Weitere Informationen zu AWS Data Pipeline Nutzungsichten finden Sie unter [AWS Data Pipeline](#).

Melden Sie sich für eine an AWS-Konto

Wenn Sie noch keine haben AWS-Konto, führen Sie die folgenden Schritte aus, um eine zu erstellen.

Um sich für eine anzumelden AWS-Konto

1. Öffnen Sie <https://portal.aws.amazon.com/billing/die Anmeldung.>
2. Folgen Sie den Online-Anweisungen.

Während der Anmeldung erhalten Sie einen Telefonanruf oder eine Textnachricht und müssen einen Verifizierungscode über die Telefontasten eingeben.

Wenn Sie sich für eine anmelden AWS-Konto, Root-Benutzer des AWS-Kontos wird eine erstellt. Der Root-Benutzer hat Zugriff auf alle AWS-Services und Ressourcen des Kontos. Als bewährte Sicherheitsmethode weisen Sie einem Administratorbenutzer Administratorzugriff zu und verwenden Sie nur den Root-Benutzer, um [Aufgaben auszuführen, die Root-Benutzerzugriff erfordern.](#)

AWS sendet Ihnen nach Abschluss des Anmeldevorgangs eine Bestätigungs-E-Mail. Du kannst jederzeit deine aktuellen Kontoaktivitäten einsehen und dein Konto verwalten, indem du zu <https://aws.amazon.com/gehst> und Mein Konto auswählst.

Erstellen eines Benutzers mit Administratorzugriff

Nachdem Sie sich für einen angemeldet haben AWS-Konto, sichern Sie Ihren Root-Benutzer des AWS-Kontos AWS IAM Identity Center, aktivieren und erstellen Sie einen Administratorbenutzer, sodass Sie den Root-Benutzer nicht für alltägliche Aufgaben verwenden.

Sichern Sie Ihre Root-Benutzer des AWS-Kontos

1. Melden Sie sich [AWS-Managementkonsole](#) als Kontoinhaber an, indem Sie Root-Benutzer auswählen und Ihre AWS-Konto E-Mail-Adresse eingeben. Geben Sie auf der nächsten Seite Ihr Passwort ein.

Hilfe bei der Anmeldung mit dem Root-Benutzer finden Sie unter [Anmelden als Root-Benutzer](#) im AWS-Anmeldung Benutzerhandbuch zu.

2. Aktivieren Sie die Multi-Faktor-Authentifizierung (MFA) für den Root-Benutzer.

Anweisungen finden Sie unter [Aktivieren eines virtuellen MFA-Geräts für Ihren AWS-Konto Root-Benutzer \(Konsole\)](#) im IAM-Benutzerhandbuch.

Erstellen eines Benutzers mit Administratorzugriff

1. Aktivieren Sie das IAM Identity Center.

Anweisungen finden Sie unter [Aktivieren AWS IAM Identity Center](#) im AWS IAM Identity Center Benutzerhandbuch.

2. Gewähren Sie einem Administratorbenutzer im IAM Identity Center Benutzerzugriff.

Ein Tutorial zur Verwendung von IAM-Identity-Center-Verzeichnis als Identitätsquelle finden Sie IAM-Identity-Center-Verzeichnis im Benutzerhandbuch unter [Benutzerzugriff mit der Standardeinstellung konfigurieren](#). AWS IAM Identity Center

Anmelden als Administratorbenutzer

- Um sich mit Ihrem IAM-Identity-Center-Benutzer anzumelden, verwenden Sie die Anmelde-URL, die an Ihre E-Mail-Adresse gesendet wurde, als Sie den IAM-Identity-Center-Benutzer erstellt haben.

Hilfe bei der Anmeldung mit einem IAM Identity Center-Benutzer finden Sie [im AWS-Anmeldung Benutzerhandbuch unter Anmeldung beim AWS Access-Portal](#).

Weiteren Benutzern Zugriff zuweisen

- Erstellen Sie im IAM-Identity-Center einen Berechtigungssatz, der den bewährten Vorgehensweisen für die Anwendung von geringsten Berechtigungen folgt.

Anweisungen hierzu finden Sie unter [Berechtigungssatz erstellen](#) im AWS IAM Identity Center Benutzerhandbuch.

- Weisen Sie Benutzer einer Gruppe zu und weisen Sie der Gruppe dann Single Sign-On-Zugriff zu.

Eine genaue Anleitung finden Sie unter [Gruppen hinzufügen](#) im AWS IAM Identity Center Benutzerhandbuch.

Erstellen Sie IAM-Rollen für AWS Data Pipeline und Pipeline-Ressourcen

AWS Data Pipeline erfordert IAM-Rollen, die die Berechtigungen für die Ausführung von Aktionen und den Zugriff auf AWS Ressourcen festlegen. Die Pipeline-Rolle bestimmt die Berechtigungen, AWS Data Pipeline über die eine Person verfügt, und eine Ressourcenrolle bestimmt die Berechtigungen, über die Anwendungen verfügen, die auf Pipeline-Ressourcen, wie z. B. EC2 Instanzen, ausgeführt werden. Sie geben diese Rollen an, wenn Sie eine Pipeline erstellen. Auch wenn Sie keine benutzerdefinierte Rolle angeben und die Standardrollen `DataPipelineDefaultRole` verwenden `DataPipelineDefaultResourceRole`, müssen Sie zuerst die Rollen erstellen und Berechtigungsrichtlinien anhängen. Weitere Informationen finden Sie unter [IAM-Rollen für AWS Data Pipeline](#).

Erlauben Sie IAM-Prinzipalen (Benutzern und Gruppen), die erforderlichen Aktionen auszuführen

Um mit einer Pipeline arbeiten zu können, muss ein IAM-Prinzipal (ein Benutzer oder eine Gruppe) in Ihrem Konto die erforderlichen [AWS Data Pipeline Aktionen und Aktionen](#) für andere Dienste ausführen dürfen, wie in Ihrer Pipeline definiert.

Um die Berechtigungen zu vereinfachen, steht Ihnen die AWSDataPipeline_FullAccessverwaltete Richtlinie zur Verfügung, die Sie an IAM-Prinzipale anhängen können. Diese verwaltete Richtlinie ermöglicht es dem Prinzipal, alle Aktionen auszuführen, die ein Benutzer benötigt, sowie die `iam:PassRole` Aktion für die Standardrollen, die verwendet werden AWS Data Pipeline , wenn keine benutzerdefinierte Rolle angegeben ist.

Es wird dringend empfohlen, diese verwaltete Richtlinie sorgfältig zu prüfen und die Berechtigungen nur auf diejenigen zu beschränken, die Ihre Benutzer benötigen. Verwenden Sie bei Bedarf diese Richtlinie als Ausgangspunkt und entfernen Sie dann die Berechtigungen, um eine restriktivere Inline-Berechtigungsrichtlinie zu erstellen, die Sie an IAM-Prinzipale anhängen können. Weitere Informationen und Beispiele für Berechtigungsrichtlinien finden Sie unter [Beispielrichtlinien für AWS Data Pipeline](#)

Eine Richtlinienanweisung, die dem folgenden Beispiel ähnelt, muss in einer Richtlinie enthalten sein, die einem beliebigen IAM-Prinzipal zugeordnet ist, der die Pipeline verwendet. Diese Anweisung ermöglicht es dem IAM-Prinzipal, die `PassRole` Aktion für die Rollen auszuführen, die eine Pipeline verwendet. Wenn Sie keine Standardrollen verwenden, ersetzen Sie `MyPipelineRole` und `MyResourceRole` durch die benutzerdefinierten Rollen, die Sie erstellen.

JSON

```
{  
    "Version": "2012-10-17",  
    "Statement": [  
        {  
            "Action": "iam:PassRole",  
            "Effect": "Allow",  
            "Resource": [  
                "arn:aws:iam::*:role/MyPipelineRole",  
                "arn:aws:iam::*:role/MyResourceRole"  
            ]  
    ]}
```

```
}
```

Das folgende Verfahren zeigt, wie Sie eine IAM-Gruppe erstellen, die AWSDataPipeline_FullAccessverwaltete Richtlinie an die Gruppe anhängen und dann Benutzer zur Gruppe hinzufügen. Sie können dieses Verfahren für jede Inline-Richtlinie verwenden

Um eine Benutzergruppe zu erstellen **DataPipelineDevelopers** und die AWSDataPipeline_FullAccessRichtlinie anzuhängen

1. Öffnen Sie unter <https://console.aws.amazon.com/iam/> die IAM-Konsole.
2. Wechseln Sie im Navigationsbereich zu Groups und klicken Sie auf Create New Group.
3. Geben Sie beispielsweise einen Gruppennamen ein und wählen Sie dann Next Step aus.
DataPipelineDevelopers
4. Geben Sie **AWSDataPipeline_FullAccess** Filter ein und wählen Sie ihn dann aus der Liste aus.
5. Wählen Sie Next Step (Nächster Schritt) und anschließend Create Group (Gruppe erstellen) aus.
6. Um Benutzer zur Gruppe hinzuzufügen:
 - a. Wählen Sie die Gruppe, die Sie erstellt haben, aus der Gruppenliste aus.
 - b. Wählen Sie Group Actions (Gruppenaktionen) und Add Users to Group (Benutzer zu Gruppe hinzufügen) aus.
 - c. Wählen Sie die Benutzer, die Sie hinzufügen möchten, aus der Liste aus und klicken Sie dann auf Benutzer zur Gruppe hinzufügen.

Erteilen programmgesteuerten Zugriffs

Benutzer benötigen programmgesteuerten Zugriff, wenn sie mit AWS außerhalb von interagieren möchten. AWS-Managementkonsole Die Art und Weise, wie programmatischer Zugriff gewährt wird, hängt vom Benutzertyp ab, der zugreift. AWS

Um Benutzern programmgesteuerten Zugriff zu gewähren, wählen Sie eine der folgenden Optionen.

Welcher Benutzer benötigt programmgesteuerten Zugriff?	Bis	Von
IAM	(Empfohlen) Verwenden Sie Konsolenanmeldeinformationen als temporäre Anmeldeinformationen, um programmatische Anfragen an AWS CLI AWS SDKs, oder zu signieren . AWS APIs	<p>Befolgen Sie die Anweisungen für die Schnittstelle, die Sie verwenden möchten.</p> <ul style="list-style-type: none"> Informationen zu den AWS CLI finden Sie unter Anmeldung für AWS lokale Entwicklung im AWS Command Line Interface Benutzerhandbuch. Weitere Informationen finden Sie unter Anmeldung für AWS lokale Entwicklung im Referenzhandbuch AWS SDKs und im Tools-Referenzhandbuch AWS SDKs
Mitarbeiteridentität (Benutzer, die in IAM Identity Center verwaltet werden)	Verwenden Sie temporäre Anmeldeinformationen, um programmatische Anfragen an das AWS CLI AWS SDKs, oder AWS APIs zu signieren.	<p>Befolgen Sie die Anweisungen für die Schnittstelle, die Sie verwenden möchten.</p> <ul style="list-style-type: none"> Informationen zu den AWS CLI finden Sie unter Konfiguration der AWS CLI zur Verwendung AWS IAM Identity Center im AWS Command Line Interface Benutzerhandbuch. Informationen zu AWS SDKs Tools und AWS APIs finden Sie unter IAM Identity Center-Authentifizierung im Referenzhandbuch AWS

Welcher Benutzer benötigt programmgesteuerten Zugriff?	Bis	Von
		SDKs und im Tools-Referenzhandbuch.
IAM	Verwenden Sie temporäre Anmeldeinformationen, um programmiatische Anfragen an das AWS CLI AWS SDKs, oder zu signieren. AWS APIs	Folgen Sie den Anweisungen unter Verwenden temporäre Anmeldeinformationen mit AWS Ressourcen im IAM-Benutzerhandbuch.
IAM	(Nicht empfohlen) Verwenden Sie langfristige Anmeldeinformationen, um programmiatische Anfragen an das AWS CLI AWS SDKs, oder zu signieren. AWS APIs	<p>Befolgen Sie die Anweisungen für die Schnittstelle, die Sie verwenden möchten.</p> <ul style="list-style-type: none"> Informationen dazu AWS CLI finden Sie unter Authentifizierung mithilfe von IAM-Benutzeranmeldenformationen im AWS Command Line Interface Benutzerhandbuch. Informationen zu AWS SDKs und Tools finden Sie unter Authentifizieren mit langfristigen Anmeldeinformationen im Referenzhandbuch AWS SDKs und im Tools-Referenzhandbuch. Weitere Informationen finden Sie unter Verwaltung von Zugriffsschlüsseln für IAM-Benutzer im IAM-Benutzerhandbuch. AWS APIs

Erste Schritte mit AWS Data Pipeline

AWS Data Pipeline hilft Ihnen dabei, wiederkehrende Datenverarbeitungs-Workloads zuverlässig und kostengünstig zu sequenzieren, zu planen, auszuführen und zu verwalten. Dieser Service erleichtert Ihnen das Entwerfen extract-transform-load (ETL) -Aktivitäten mithilfe strukturierter und unstrukturierter Daten, sowohl vor Ort als auch in der Cloud, auf der Grundlage Ihrer Geschäftslogik.

Zur Verwendung AWS Data Pipeline erstellen Sie eine Pipeline-Definition, die die Geschäftslogik für Ihre Datenverarbeitung festlegt. Eine typische Pipeline-Definition besteht aus [Aktivitäten](#), die die auszuführende Arbeit definieren, und [Datenknoten](#), die den Ort und die Art der Eingabe- und Ausgabedaten definieren.

In diesem Tutorial führen Sie ein Shell-Befehlsskript aus, das die Anzahl der GET-Anforderungen in Apache-Webserverprotokollen zählt. Diese Pipeline läuft eine Stunde lang alle 15 Minuten und schreibt bei jeder Iteration die Ausgabe in Amazon S3.

Voraussetzungen

Bevor Sie beginnen, führen Sie die Aufgaben in [Einrichtung für AWS Data Pipeline](#) durch.

Pipeline-Objekte

Die Pipeline verwendet die folgenden Objekte:

[ShellCommandActivity](#)

Liest die Eingabeprotokolldatei und zählt die Anzahl an Fehlern.

[S3 DataNode](#) (Eingabe)

Der S3-Bucket, der die Eingabeprotokolldatei enthält.

[S3 DataNode](#) (Ausgabe)

Der S3-Bucket für die Ausgabe.

[Ec2Resource](#)

Die Rechenressource, die zur Ausführung der Aktivität AWS Data Pipeline verwendet wird.

Beachten Sie, dass Sie, wenn Sie über eine große Menge an Protokolldateidaten verfügen, Ihre Pipeline so konfigurieren können, dass anstelle einer EC2 Instanz ein EMR-Cluster zur Verarbeitung der Dateien verwendet wird.

Plan

Legt fest, dass die Aktivität alle 15 Minuten eine Stunde lang ausgeführt wird.

Aufgaben

- [Erstellen Sie die Pipeline](#)
- [Überwachen der ausgeföhrten Pipeline](#)
- [Anzeigen der Ausgabe](#)
- [Löschen der Pipeline](#)

Erstellen Sie die Pipeline

Der schnellste Weg, damit zu beginnen, AWS Data Pipeline ist die Verwendung einer Pipeline-Definition, die als Vorlage bezeichnet wird.

So erstellen Sie die Pipeline

1. Öffnen Sie die AWS Data Pipeline Konsole unter. <https://console.aws.amazon.com/datapipeline/>
2. Wählen Sie auf der Navigationsleiste eine Region aus. Sie können unabhängig von Ihrem Standort jede verfügbare Region auswählen. Viele AWS-Ressourcen sind regionsspezifisch, AWS Data Pipeline ermöglichen es Ihnen jedoch, Ressourcen zu verwenden, die sich in einer anderen Region als der Pipeline befinden.
3. Der erste Bildschirm, der angezeigt wird, hängt davon ab, ob Sie in der aktuellen Region eine Pipeline erstellt haben.
 - a. Wenn Sie in dieser Region keine Pipeline erstellt haben, zeigt die Konsole einen Einführungsbildschirm an. Wählen Sie Get started now.
 - b. Wenn Sie in dieser Region bereits eine Pipeline erstellt haben, zeigt die Konsole eine Seite an, auf der Ihre Pipelines für die Region aufgeführt sind. Wählen Sie Create new pipeline (Neue Pipeline erstellen) aus.
4. Geben Sie im Feld Name einen Namen für Ihre Pipeline ein.
5. (Optional) Geben Sie unter Beschreibung eine Beschreibung für Ihre Pipeline ein.
6. Wählen Sie unter Quelle die Option Mit einer Vorlage erstellen und dann die folgende Vorlage aus: Erste Schritte mit ShellCommandActivity.

7. Nach der Auswahl der Vorlage öffnet sich der Abschnitt Parameters. Behalten Sie dort die Standardwerte für S3 input folder und Shell command to run bei. Klicken Sie neben S3 output folder auf das Ordnersymbol, wählen Sie einen Ihrer Buckets oder Ordner aus und klicken Sie anschließend auf Select.
8. Behalten Sie unter Schedule die Standardwerte bei. Wenn Sie die Pipeline aktivieren, beginnen die Pipeline-Ausführungen und werden alle 15 Minuten eine Stunde lang ausgeführt.

Wenn Sie möchten, können Sie stattdessen auch die Option Run once on pipeline activation auswählen.

9. Lassen Sie unter Pipeline-Konfiguration die Protokollierung aktiviert. Wählen Sie das Ordnersymbol unter S3-Speicherort für Protokolle, wählen Sie einen Ihrer Buckets oder Ordner aus und wählen Sie dann Auswählen aus.

Wenn Sie möchten, können Sie stattdessen die Protokollierung deaktivieren.

10. Lassen Sie unter Sicherheit/Zugriff die IAM-Rollen auf Standard eingestellt.
11. Klicken Sie auf Activate.

Wenn Sie möchten, können Sie „In Architect bearbeiten“ wählen, um diese Pipeline zu ändern. Sie können beispielsweise Vorbedingungen hinzufügen.

Überwachen der ausgeführten Pipeline

Nachdem Sie Ihre Pipeline aktiviert haben, können Sie auf die Seite Execution details gehen, wo Sie den Fortschritt Ihrer Pipeline überwachen können.

So überwachen Sie den Fortschritt Ihrer Pipeline

1. Klicken Sie auf Update oder drücken Sie F5, um den angezeigten Status zu aktualisieren.

 Tip

Wenn keine Ausführungen aufgelistet sind, stellen Sie sicher, dass Start (in UTC) und End (in UTC) die geplante Start- und Endzeit Ihrer Pipeline abdecken, und klicken Sie dann auf Update.

2. Wenn der Status jedes Objekt in der Pipeline FINISHED ist, hat Ihre Pipeline die geplanten Tasks erfolgreich fertiggestellt.

3. Wenn Ihre Pipeline nicht erfolgreich abgeschlossen wurde, überprüfen Sie Ihre Pipeline-Einstellungen auf Probleme. Weitere Informationen zur Fehlerbehebung bei fehlgeschlagenen oder unvollständigen Instance-Ausführungen Ihrer Pipeline finden Sie unter [Beheben typischer Probleme](#).

Anzeigen der Ausgabe

Öffnen Sie die Amazon S3 S3-Konsole und navigieren Sie zu Ihrem Bucket. Wenn Sie Ihre Pipeline alle 15 Minuten eine Stunde lang ausgeführt haben, sehen Sie vier Unterordner mit Zeitstempeln. Jeder Unterordner enthält die Ausgabe in einer Datei mit dem Namen `output.txt`. Da wir das Skript jedes Mal auf derselben Eingabedatei ausgeführt haben, sind die Ausgabedateien identisch.

Löschen der Pipeline

Um keine Gebühren mehr anfallen zu lassen, löschen Sie Ihre Pipeline. Beim Löschen Ihrer Pipeline werden die Pipeline-Definition und alle zugehörigen Objekte gelöscht.

Um Ihre Pipeline zu löschen

1. Wählen Sie auf der Seite Pipelines auflisten Ihre Pipeline aus.
2. Klicken Sie auf Aktionen und wählen Sie dann Löschen aus.
3. Wenn Sie zur Bestätigung aufgefordert werden, wählen Sie Delete (Löschen).

Wenn Sie mit der Ausgabe dieses Tutorials fertig sind, löschen Sie die Ausgabeordner aus Ihrem Amazon S3 S3-Bucket.

Arbeiten mit Pipelines

Sie können Pipelines mithilfe der Befehlszeilenschnittstelle (CLI) oder AWS des SDK verwalten, erstellen und ändern. Die folgenden Abschnitte beschreiben die grundlegenden AWS Data Pipeline - Konzepte und veranschaulichen, wie Sie mit Pipelines arbeiten.

Important

Bevor Sie beginnen, sehen Sie sich [Einrichtung für AWS Data Pipeline](#) an.

Inhalt

- [Eine Pipeline erstellen](#)
- [Anzeigen Ihrer Pipelines](#)
- [Bearbeiten Ihrer Pipeline](#)
- [Klonen Ihrer Pipeline](#)
- [Tagging Ihrer Pipeline](#)
- [Deaktivieren Ihrer Pipeline](#)
- [Löschen Ihrer Pipeline](#)
- [Staging von Daten und Tabellen mit Pipeline-Aktivitäten](#)
- [Verwenden einer Pipeline mit Ressourcen in mehreren Regionen](#)
- [Cascading-Ausfälle und erneute Ausführungen](#)
- [Syntax der Pipeline-Definitionsdatei](#)
- [Arbeiten mit der API](#)

Eine Pipeline erstellen

AWS Data Pipeline bietet Ihnen mehrere Möglichkeiten, Pipelines zu erstellen:

- Verwenden Sie die AWS Command Line Interface (CLI) mit einer Vorlage, die Ihnen zur Verfügung gestellt wird. Weitere Informationen finden Sie unter [Erstellen Sie mit der CLI eine Pipeline aus Data Pipeline-Vorlagen](#).

- Verwenden Sie die AWS Command Line Interface (CLI) mit einer Pipeline-Definitionsdatei im JSON-Format.
- Sie können ein AWS SDK mit einer sprachspezifischen API verwenden. Weitere Informationen finden Sie unter [Arbeiten mit der API](#).

Erstellen Sie mit der CLI eine Pipeline aus Data Pipeline-Vorlagen

Data Pipeline bietet mehrere vorkonfigurierte Pipeline-Definitionen, sogenannte Vorlagen. Sie können Vorlagen verwenden, um AWS Data Pipeline schnell loszulegen. Diese Vorlagen sind in einem öffentlichen Bucket am Amazon S3 S3-Standort verfügbar:`s3://datapipeline-us-east-1/templates/`. Diese vordefinierten Vorlagen wurden für bestimmte Anwendungsfälle erstellt und können zur Erstellung von Pipelines verwendet werden. Sie können sie verwenden `aws s3 ls --recursive "s3://datapipeline-us-east-1/templates/"`, um alle verfügbaren Vorlagen aufzulisten.

Erstellen Sie mit der CLI eine Pipeline aus einer Vorlage

Angenommen, Sie möchten eine Pipeline erstellen, die eine DynamoDB-Tabelle nach Amazon S3 exportiert. Die in diesem Fall zu verwendende Vorlage finden Sie unter: `s3://datapipeline-us-east-1/templates/DynamoDB Templates/Export DynamoDB table to S3.json`

Um die JSON-Vorlage herunterzuladen und eine Pipeline mit der CLI zu erstellen

1. Laden Sie die Vorlage mit der `aws s3 cp` CLI oder Curl herunter. Zum Beispiel:

```
aws s3 cp "s3://datapipeline-us-east-1/templates/DynamoDB Templates/Export DynamoDB table to S3.json" <destination directory>
```

2. Nehmen Sie nach Bedarf Änderungen an der heruntergeladenen Vorlage vor. Um beispielsweise die neueste EMR-Release-Version zu verwenden, ändern Sie das `releaseLabel` Feld im `EmrClusterForBackup` Objekt, ändern Sie die Master- und Core-Instance-Typen und ändern Sie die Standardwerte der Parameter in der Vorlage.

3. Erstellen Sie eine Pipeline mit der `create-pipeline` CLI. Zum Beispiel:

```
aws datapipeline create-pipeline --name my-ddb-backup-pipeline --unique-id my-ddb-backup-pipeline --region ap-northeast-1
```

4. Notieren Sie sich die erstellte Pipeline-ID.

5. Wird verwendet `put-pipeline-definition`, um die Definition hochzuladen. Geben Sie Werte der Parameter an, deren Standardwerte Sie mithilfe der `--parameter-values` Option überschreiben möchten.

Weitere Informationen zu Vorlagen finden Sie unter [Auswahl einer Vorlage](#).

Auswahl einer Vorlage

Die folgenden Vorlagen stehen im Amazon S3 S3-Bucket zum Download zur Verfügung:`s3://datapipeline-us-east-1/templates/`.

Vorlagen

- [Erste Schritte mit ShellCommandActivity](#)
- [AWS CLI-Befehl ausführen](#)
- [DynamoDB-Tabelle nach S3 exportieren](#)
- [DynamoDB-Backup-Daten aus S3 importieren](#)
- [Job auf einem Amazon EMR-Cluster ausführen](#)
- [Vollständige Kopie der Amazon RDS-MySQL-Tabelle für Amazon S3](#)
- [Inkrementelles Kopieren der Amazon RDS-MySQL-Tabelle nach Amazon S3](#)
- [Laden Sie S3-Daten in die Amazon RDS-MySQL-Tabelle](#)
- [Vollständige Kopie der Amazon RDS-MySQL-Tabelle nach Amazon Redshift](#)
- [Inkrementelles Kopieren einer Amazon RDS-MySQL-Tabelle nach Amazon Redshift](#)
- [Daten von Amazon S3 in Amazon Redshift laden](#)

Erste Schritte mit ShellCommandActivity

Die ShellCommandActivity Vorlage „Erste Schritte verwenden“ führt ein Shell-Befehlsskript aus, um die Anzahl der GET-Anfragen in einer Protokolldatei zu zählen. Die Ausgabe wird bei jedem geplanten Lauf der Pipeline an einen Amazon S3 S3-Speicherort mit Zeitstempel geschrieben.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- ShellCommandActivity
- S3 InputNode
- S3 OutputNode

- Ec2Resource

AWS CLI-Befehl ausführen

Diese Vorlage führt in geplanten Intervallen einen benutzerdefinierten AWS CLI Befehl aus.

DynamoDB-Tabelle nach S3 exportieren

Die Vorlage DynamoDB-Tabelle nach S3 exportieren plant, dass ein Amazon EMR-Cluster Daten aus einer DynamoDB-Tabelle in einen Amazon S3 S3-Bucket exportiert. Diese Vorlage verwendet einen Amazon EMR-Cluster, dessen Größe proportional zum Wert des für die DynamoDB-Tabelle verfügbaren Durchsatzes ist. Sie können die Anzahl einer Tabelle zwar erhöhen IOPs , dies kann jedoch zu zusätzlichen Kosten beim Import und Export führen. Bisher wurde für den Export ein verwendet HiveActivity , jetzt aber nativ MapReduce.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- [EmrActivity](#)
- [EmrCluster](#)
- [Dynamo-Knoten DBData](#)
- [S3 DataNode](#)

DynamoDB-Backup-Daten aus S3 importieren

Die Vorlage DynamoDB-Backup-Daten aus S3 importieren plant, dass ein Amazon EMR-Cluster ein zuvor erstelltes DynamoDB-Backup in Amazon S3 in eine DynamoDB-Tabelle lädt. Bestehende Elemente in der DynamoDB-Tabelle werden mit denen aus den Sicherungsdaten aktualisiert, und der Tabelle werden neue Elemente hinzugefügt. Diese Vorlage verwendet einen Amazon EMR-Cluster, dessen Größe proportional zum Wert des für die DynamoDB-Tabelle verfügbaren Durchsatzes ist. Sie können die Anzahl einer Tabelle zwar erhöhen IOPs , dies kann jedoch zu zusätzlichen Kosten beim Import und Export führen. Bisher wurde für den Import ein verwendet, HiveActivity jetzt wird nativ MapReduce verwendet.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- [EmrActivity](#)
- [EmrCluster](#)

- [Dynamo-Knoten DBData](#)
- [S3 DataNode](#)
- [S3 PrefixNotEmpty](#)

Job auf einem Amazon EMR-Cluster ausführen

Die Vorlage „Job auf einem Elastic MapReduce Cluster ausführen“ startet einen Amazon EMR-Cluster auf der Grundlage der angegebenen Parameter und startet die Ausführung von Schritten auf der Grundlage des angegebenen Zeitplans. Sobald der Auftrag abgeschlossen ist, wird der EMR-Cluster beendet. Optionale Bootstrap-Aktionen können angegeben werden, um zusätzliche Software zu installieren oder die Konfiguration der Anwendung im Cluster zu ändern.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- [EmrActivity](#)
- [EmrCluster](#)

Vollständige Kopie der Amazon RDS-MySQL-Tabelle für Amazon S3

Die Vorlage Vollständige Kopie der RDS-MySQL-Tabelle in S3 kopiert eine gesamte Amazon RDS-MySQL-Tabelle und speichert die Ausgabe an einem Amazon S3 S3-Speicherort. Die Ausgabe wird als CSV-Datei in einem Unterordner mit Zeitstempel unter dem angegebenen Amazon S3 S3-Speicherort gespeichert.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- [CopyActivity](#)
- [Ec2Resource](#)
- [SqlDataNode](#)
- [S3 DataNode](#)

Inkrementelles Kopieren der Amazon RDS-MySQL-Tabelle nach Amazon S3

Die Vorlage Incremental Copy of RDS MySQL Table to S3 erstellt eine inkrementelle Kopie der Daten aus einer Amazon RDS-MySQL-Tabelle und speichert die Ausgabe an einem Amazon S3 S3-Speicherort. Die Amazon RDS-MySQL-Tabelle muss die Spalte Zuletzt geändert haben.

Diese Vorlage kopiert alle Änderungen, die ab dem geplanten Startzeitpunkt zwischen festgelegten Intervallen an der Tabelle vorgenommen werden. Der Zeitplantyp ist Zeitreihe. Wenn also eine Kopie für eine bestimmte Stunde geplant war, werden die Tabellenzeilen AWS Data Pipeline kopiert, deren Zeitstempel „Letzte Änderung“ in diese Stunde fällt. Physische Löschtätigkeiten an der Tabelle werden nicht kopiert. Die Ausgabe wird bei jeder geplanten Ausführung in einen Unterordner mit Zeitstempel unter dem Amazon S3 S3-Speicherort geschrieben.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- [CopyActivity](#)
- [Ec2Resource](#)
- [SqlDataNode](#)
- [S3 DataNode](#)

Laden Sie S3-Daten in die Amazon RDS-MySQL-Tabelle

Die Vorlage S3-Daten in eine RDS-MySQL-Tabelle laden plant, dass eine EC2 Amazon-Instance die CSV-Datei aus dem unten angegebenen Amazon S3-Dateipfad in eine Amazon RDS-MySQL-Tabelle kopiert. Die CSV-Datei sollte über keine Kopfzeile verfügen. Die Vorlage aktualisiert bestehende Einträge in der Amazon RDS-MySQL-Tabelle mit denen in den Amazon S3 S3-Daten und fügt neue Einträge aus den Amazon S3-Daten zur Amazon RDS-MySQL-Tabelle hinzu. Sie können die Daten in eine vorhandene Tabelle laden oder eine SQL-Abfrage zum Erstellen einer neuen Tabelle bereitstellen.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- [CopyActivity](#)
- [Ec2Resource](#)
- [SqlDataNode](#)
- [S3 DataNode](#)

Vorlagen von Amazon RDS zu Amazon Redshift

Die folgenden beiden Vorlagen kopieren Tabellen von Amazon RDS MySQL nach Amazon Redshift mithilfe eines Übersetzungsskripts, das eine Amazon Redshift Redshift-Tabelle unter Verwendung des Quelltabellenschemas mit den folgenden Einschränkungen erstellt:

- Wenn kein Verteilungsschlüssel angegeben ist, wird der erste Primärschlüssel aus der Amazon RDS-Tabelle als Verteilungsschlüssel festgelegt.
- Sie können eine Spalte, die in einer Amazon RDS-MySQL-Tabelle vorhanden ist, nicht überspringen, wenn Sie eine Kopie nach Amazon Redshift erstellen.
- (Optional) Sie können eine Amazon RDS-Datentyp-Zuordnung von MySQL zu Amazon Redshift als einen der Parameter in der Vorlage angeben. Wenn dies angegeben ist, verwendet das Skript dies, um die Amazon Redshift Redshift-Tabelle zu erstellen.

Wenn der `Overwrite_Existing` Amazon Redshift Redshift-Einfügemodus verwendet wird:

- Wenn kein Verteilungsschlüssel angegeben wird, wird ein Primärschlüssel in der Amazon RDS-MySQL-Tabelle verwendet.
- Wenn zusammengesetzte Primärschlüssel für die Tabelle vorhanden sind, wird der erste davon als Verteilungsschlüssel verwendet, sofern kein Verteilungsschlüssel bereitgestellt wird. Nur der erste zusammengesetzte Schlüssel ist als Primärschlüssel in der Amazon Redshift Redshift-Tabelle festgelegt.
- Wenn kein Verteilungsschlüssel bereitgestellt wird und es keinen Primärschlüssel in der Amazon RDS-MySQL-Tabelle gibt, schlägt der Kopievorgang fehl.

Weitere Informationen zu Amazon Redshift finden Sie in den folgenden Themen:

- [Amazon-Redshift-Cluster](#)
- [Amazon Redshift KOPIEREN](#)
- [Verteilungsstile](#) und [DISTKEY-Beispiele](#)
- [Sortierschlüssel](#)

In der folgenden Tabelle wird beschrieben, wie das Skript die Datentypen umwandelt:

Datentypübersetzungen zwischen MySQL und Amazon Redshift

MySQL-Datentyp	Amazon Redshift Redshift-Datentyp	Hinweise
TINYINT,	SMALLINT	
TINYINT (Größe)		MySQL: -128 bis 127. Die maximale Anzahl von Ziffern

MySQL-Datentyp	Amazon Redshift Redshift-Datentyp	Hinweise
		kann in Klammern angegeben werden. Amazon Redshift: INT2. 2-Byte-Ganzzahl mit Vorzeichen
TINYINT UNSIGNED, TINYINT (Größe) UNSIGNED	SMALLINT	MySQL: 0 bis 255 UNSIGNED. Die maximale Anzahl von Ziffern kann in Klammern angegeben werden. Amazon Redshift: INT2. 2-Byte-Ganzzahl mit Vorzeichen
SMALLINT, SMALLINT (Größe)	SMALLINT	MySQL: -32768 bis 32767 normal. Die maximale Anzahl von Ziffern kann in Klammern angegeben werden. Amazon Redshift: INT2. 2-Byte-Ganzzahl mit Vorzeichen
SMALLINT UNSIGNED, SMALLINT(Größe) UNSIGNED,	INTEGER	MySQL: 0 bis 65535 UNSIGNED*. Die maximale Anzahl von Ziffern kann in Klammern angegeben werden. Amazon Redshift: INT4. 4-Byte-Ganzzahl mit Vorzeichen

MySQL-Datentyp	Amazon Redshift Redshift-Datentyp	Hinweise
MEDIUMINT, MEDIUMINT (Größe)	INTEGER	MySQL: 388608 bis 8388607. Die maximale Anzahl von Ziffern kann in Klammern angegeben werden. Amazon Redshift: INT4. 4-Byte-Ganzzahl mit Vorzeichen
MEDIUMINT UNSIGNED, MEDIUMINT (Größe) UNSIGNED	INTEGER	MySQL: 0 bis 16777215. Die maximale Anzahl von Ziffern kann in Klammern angegeben werden. Amazon Redshift: INT4. 4-Byte-Ganzzahl mit Vorzeichen
INT, INT(Größe)	INTEGER	MySQL: 147483648 bis 2147483647 Amazon Redshift: INT4. 4-Byte-Ganzzahl mit Vorzeichen
INT UNSIGNED, INT(Größe) UNSIGNED	BIGINT	MySQL: 0 bis 4294967295 Amazon Redshift: INT8. 8-Byte-Ganzzahl mit Vorzeichen
BIGINT BIGINT(Größe)	BIGINT	Amazon Redshift: INT8. 8-Byte-Ganzzahl mit Vorzeichen
BIGINT UNSIGNED BIGINT(Größe) UNSIGNED	VARCHAR(20*4)	MySQL: 0 bis 184467440 73709551615 Amazon Redshift: Kein natives Äquivalent, also Verwendung eines Char-Arrays.

MySQL-Datentyp	Amazon Redshift Redshift-Datentyp	Hinweise
FLOAT FLOAT(Größe,d) FLOAT(Größe,d) UNSIGNED	REAL	<p>Die maximale Anzahl von Ziffern kann im Größenparameter angegeben werden.</p> <p>Die maximale Anzahl von Ziffern rechts neben dem Dezimalzeichen wird im d-Parameter angegeben.</p> <p>Amazon Redshift: FLOAT4</p>
DOUBLE(Größe,d)	DOUBLE PRECISION	<p>Die maximale Anzahl von Ziffern kann im Größenparameter angegeben werden.</p> <p>Die maximale Anzahl von Ziffern rechts neben dem Dezimalzeichen wird im d-Parameter angegeben.</p> <p>Amazon Redshift: FLOAT8</p>
DECIMAL(Größe,d)	DECIMAL(Größe,d)	<p>Als Zeichenfolge gespeicherte DOUBLE, ermöglicht ein festes Dezimalzeichen. Die maximale Anzahl von Ziffern kann im Größenparameter angegeben werden. Die maximale Anzahl von Ziffern rechts neben dem Dezimalzeichen wird im d-Parameter angegeben.</p> <p>Amazon Redshift: Kein natives Äquivalent.</p>

MySQL-Datentyp	Amazon Redshift Redshift-Datentyp	Hinweise
CHAR(Größe)	VARCHAR(Größe*4)	<p>Enthält eine Zeichenfolge fester Länge, die aus Buchstaben, Ziffern und Sonderzeichen bestehen kann. Die feste Größe wird als Parameter in Klammern angegeben. Kann bis zu 255 Zeichen speichern.</p> <p>Rechts aufgefüllt mit Leerzeichen.</p> <p>Amazon Redshift: Der CHAR-Datentyp unterstützt kein Multibyte-Zeichen, daher wird VARCHAR verwendet.</p> <p>Die maximale Anzahl von Byte pro Zeichen ist entsprechend 4 RFC3629, was die Zeichentabelle auf U+10FFFF beschränkt.</p>
VARCHAR(Größe)	VARCHAR(Größe*4)	<p>Kann bis zu 255 Zeichen speichern.</p> <p>VARCHAR unterstützt nicht die folgenden ungültigen UTF-8-Codepunkte: 0xD800 – 0xFFFF, (Bytefolgen: ED A0 80 – ED BF BF), 0xFDD0 – 0xFDEF, 0xFFFFE und 0xFFFFF, (Bytefolgen: EF B7 90 – EF B7 AF, EF BF BE und EF BF BF)</p>

MySQL-Datentyp	Amazon Redshift Redshift-Datentyp	Hinweise
TINYTEXT	VARCHAR(255*4)	Enthält eine Zeichenfolge mit einer maximalen Länge von 255 Zeichen.
TEXT	VARCHAR(max)	Enthält eine Zeichenfolge mit einer maximalen Länge von 65.535 Zeichen.
MEDIUMTEXT	VARCHAR(max)	0 bis 16.777.215 Zeichen
LONGTEXT	VARCHAR(max)	0 bis 4.294.967.295 Zeichen
BOOLEAN BOOL TINYINT(1)	BOOLEAN	MySQL: Diese Typen sind Synonyme für TINYINT (1) . Der Wert Null wird als „false“ angesehen. Werte ungleich Null werden als „true“ angesehen.
BINARY[(M)]	varchar(255)	M ist 0 bis 255 Byte, FIXED
VARBINARY(M)	VARCHAR(max)	0 bis 65.535 Byte
TINYBLOB	VARCHAR(255)	0 bis 255 Byte
BLOB	VARCHAR(max)	0 bis 65.535 Byte
MEDIUMBLOB	VARCHAR(max)	0 bis 16.777.215 Byte
LONGBLOB	VARCHAR(max)	0 bis 4.294.967.295 Byte

MySQL-Datentyp	Amazon Redshift Redshift-Datentyp	Hinweise
ENUM	VARCHAR(255*2)	Die Begrenzung gilt nicht für die Länge der Literal-Aufzählungszeichenfolge, sondern vielmehr für die in der Tabelle definierte Anzahl von Aufzählungswerten.
SET	VARCHAR(255*2)	Wie enum.
DATUM	DATUM	(JJJJ-MM-TT) „1000-01-01“ bis „9999-12-31“
TIME	VARCHAR(10*4)	(hh:mm:ss) „-838:59:59“ bis „838:59:59“
DATETIME	TIMESTAMP (ZEITSTEMPEL)	(hh:mm:ss) YYYY-MM-DD „1000-01-01 00:00:00“ bis „9999-12-31 23:59:59“
TIMESTAMP (ZEITSTEMPEL)	TIMESTAMP (ZEITSTEMPEL)	(YYYYMMDDhhmmss) 19700101000000 bis 2037+
JAHR	VARCHAR(4*4)	(YYYY) 1900 bis 2155

MySQL-Datentyp	Amazon Redshift Redshift-Datentyp	Hinweise
Spalte SERIAL	<p>ID-Generierung / Dieses Attribut ist für ein OLAP-Data Warehouse nicht erforderlich, da diese Spalte kopiert wird.</p> <p>Das Schlüsselwort SERIAL wird bei der Umwandlung nicht hinzugefügt.</p>	<p>SERIAL ist tatsächlich eine Entität namens SEQUENCE. Sie existiert unabhängig von der restlichen Tabelle.</p> <p>Spalte GENERATED BY DEFAULT entspricht:</p> <pre>CREATE SEQUENCE Name; CREATE TABLE Tabelle (Spalte INTEGER NOT NULL DEFAULT nextval(Name));</pre>
Spalte BIGINT UNSIGNED NOT NULL AUTO_INCREMENT UNIQUE	<p>ID-Generierung / Dieses Attribut ist für OLAP-Data Warehouse nicht erforderlich, da diese Spalte kopiert wird.</p> <p>Das Schlüsselwort SERIAL wird bei der Umwandlung daher nicht hinzugefügt.</p>	<p>SERIAL ist tatsächlich eine Entität namens SEQUENCE. Sie existiert unabhängig von der restlichen Tabelle.</p> <p>Spalte GENERATED BY DEFAULT entspricht:</p> <pre>CREATE SEQUENCE Name; CREATE TABLE Tabelle (Spalte INTEGER NOT NULL DEFAULT nextval(Name));</pre>

MySQL-Datentyp	Amazon Redshift Redshift-Datentyp	Hinweise
ZEROFILL	Das Schlüsselwort ZEROFILL wird bei der Umwandlung nicht hinzugefügt.	<p>INT UNSIGNED ZEROFILL NOT NULL</p> <p>ZEROFILL füllt den angezeigten Wert des Feldes bis zur Anzeigenbreite in der Spaltendefinition mit Nullen auf. Werte, die länger als die Anzeigenbreite sind, werden nicht abgeschnitten. Beachten Sie, dass die Syntax von ZEROFILL auch UNSIGNED impliziert.</p>

Vollständige Kopie der Amazon RDS-MySQL-Tabelle nach Amazon Redshift

Die Vorlage Vollständige Kopie der Amazon RDS-MySQL-Tabelle in Amazon Redshift kopiert die gesamte Amazon RDS-MySQL-Tabelle in eine Amazon Redshift Redshift-Tabelle, indem Daten in einem Amazon S3 S3-Ordner bereitgestellt werden. Der Amazon S3 S3-Staging-Ordner muss sich in derselben Region wie der Amazon Redshift Redshift-Cluster befinden. Eine Amazon Redshift Redshift-Tabelle wird mit demselben Schema wie die Amazon RDS-MySQL-Quelltabelle erstellt, sofern sie noch nicht existiert. Bitte geben Sie alle Datentyp-Überschreibungen von Amazon RDS MySQL zu Amazon Redshift an, die Sie bei der Erstellung der Amazon Redshift Redshift-Tabelle anwenden möchten.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- [CopyActivity](#)
- [RedshiftCopyActivity](#)
- [S3 DataNode](#)
- [SqlDataNode](#)
- [RedshiftDataNode](#)
- [RedshiftDatabase](#)

Inkrementelles Kopieren einer Amazon RDS-MySQL-Tabelle nach Amazon Redshift

Die Vorlage Inkrementelle Kopie der Amazon RDS-MySQL-Tabelle nach Amazon Redshift kopiert Daten aus einer Amazon RDS-MySQL-Tabelle in eine Amazon Redshift Redshift-Tabelle, indem Daten in einem Amazon S3 S3-Ordner bereitgestellt werden.

Der Amazon S3 S3-Staging-Ordner muss sich in derselben Region wie der Amazon Redshift Redshift-Cluster befinden.

AWS Data Pipeline verwendet ein Übersetzungsskript, um eine Amazon Redshift Redshift-Tabelle mit demselben Schema wie die Amazon RDS-MySQL-Quelltable zu erstellen, falls diese noch nicht existiert. Sie müssen alle Datentyp-Überschreibungen von Amazon RDS MySQL bis Amazon Redshift angeben, die Sie bei der Erstellung der Amazon Redshift Redshift-Tabelle anwenden möchten.

Diese Vorlage kopiert Änderungen, die zwischen geplanten Intervallen, beginnend mit der geplanten Startzeit, an der Amazon RDS-MySQL-Tabelle vorgenommen wurden. Physische Löschen in der Amazon RDS-MySQL-Tabelle werden nicht kopiert. Sie müssen den Namen der Spalte angeben, in der der Zeitpunkt der letzten Änderung gespeichert wird.

Wenn Sie die Standardvorlage verwenden, um Pipelines für inkrementelle Amazon RDS-Kopien zu erstellen, RDSToS3CopyActivity wird eine Aktivität mit dem Standardnamen erstellt. Sie können sie umbenennen.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- [CopyActivity](#)
- [RedshiftCopyActivity](#)
- [S3 DataNode](#)
- [SqlDataNode](#)
- [RedshiftDataNode](#)
- [RedshiftDatabase](#)

Daten von Amazon S3 in Amazon Redshift laden

Die Vorlage Daten aus S3 in Redshift laden kopiert Daten aus einem Amazon S3 S3-Ordner in eine Amazon Redshift Redshift-Tabelle. Sie können die Daten in eine vorhandene Tabelle laden oder eine SQL-Abfrage zum Erstellen der Tabelle bereitstellen.

Die Daten werden basierend auf den Amazon Redshift COPY Redshift-Optionen kopiert. Die Amazon Redshift Redshift-Tabelle muss dasselbe Schema wie die Daten in Amazon S3 haben. COPYOptionen finden Sie unter [COPY](#) im Amazon Redshift Database Developer Guide.

Die Vorlage verwendet die folgenden Pipeline-Objekte:

- [CopyActivity](#)
- [RedshiftCopyActivity](#)
- [S3 DataNode](#)
- [RedshiftDataNode](#)
- [RedshiftDatabase](#)
- [Ec2Resource](#)

Eine Pipeline mithilfe parametrisierter Vorlagen erstellen

Sie können eine Pipeline-Definition anhand einer parametrisierten Vorlage anpassen. Auf diese Weise können Sie eine gemeinsame Pipeline-Definition erstellen, aber unterschiedliche Parameter konfigurieren, wenn Sie die Pipeline-Definition zu einer neuen Pipeline hinzufügen.

Inhalt

- [Fügen Sie MyVariables zur Pipeline-Definition hinzu](#)
- [Definieren Sie Parameterobjekte](#)
- [Definieren von Parameterwerten](#)
- [Die Pipeline-Definition wird eingereicht](#)

Fügen Sie MyVariables zur Pipeline-Definition hinzu

Wenn Sie die Pipeline-Definitionsdatei erstellen, geben Sie Variablen mit der folgenden Syntax an: # {*myVariable*}. Der Variablen muss my vorangestellt werden. Die folgende Pipeline-Definitionsdatei, pipeline-definition.json, enthält beispielsweise die folgenden Variablen: *myShellCmd*, *myS3InputLoc*, und *myS3OutputLoc*.



In einer Pipeline-Definition sind maximal 50 Parameter zulässig.

```
{  
  "objects": [  
    {  
      "id": "ShellCommandActivityObj",  
      "input": {  
        "ref": "S3InputLocation"  
      },  
      "name": "ShellCommandActivityObj",  
      "runsOn": {  
        "ref": "EC2ResourceObj"  
      },  
      "command": "#{myShellCmd}",  
      "output": {  
        "ref": "S3OutputLocation"  
      },  
      "type": "ShellCommandActivity",  
      "stage": "true"  
    },  
    {  
      "id": "Default",  
      "scheduleType": "CRON",  
      "failureAndRerunMode": "CASCADE",  
      "schedule": {  
        "ref": "Schedule_15mins"  
      },  
      "name": "Default",  
      "role": "DataPipelineDefaultRole",  
      "resourceRole": "DataPipelineDefaultResourceRole"  
    },  
    {  
      "id": "S3InputLocation",  
      "name": "S3InputLocation",  
      "directoryPath": "#{myS3InputLoc}",  
      "type": "S3DataNode"  
    },  
    {  
      "id": "S3OutputLocation",  
      "name": "S3OutputLocation",  
      "directoryPath": "#{myS3OutputLoc}/#{format(@scheduledStartTime, 'YYYY-MM-dd-HH-mm-ss')}",  
      "type": "S3DataNode"  
    },  
    {
```

```

    "id": "Schedule_15mins",
    "occurrences": "4",
    "name": "Every 15 minutes",
    "startAt": "FIRST_ACTIVATION_DATE_TIME",
    "type": "Schedule",
    "period": "15 Minutes"
},
{
  "terminateAfter": "20 Minutes",
  "id": "EC2ResourceObj",
  "name": "EC2ResourceObj",
  "instanceType": "t1.micro",
  "type": "Ec2Resource"
}
]
}

```

Definieren Sie Parameterobjekte

Sie können eine separate Datei mit Parameterobjekten erstellen, um die Variablen in Ihrer Pipeline-Definition zu definieren. Die folgende JSON-Datei enthält beispielsweise Parameterobjekte für die *myS3OutputLoc* Variablen *myShellCmd* und *myS3InputLoc*, und aus der obigen Beispiel-Pipeline-Definition. `parameters.json`

```

{
  "parameters": [
    {
      "id": "myShellCmd",
      "description": "Shell command to run",
      "type": "String",
      "default": "grep -rc \\"GET\\" ${INPUT1_STAGING_DIR}/* > ${OUTPUT1_STAGING_DIR}/output.txt"
    },
    {
      "id": "myS3InputLoc",
      "description": "S3 input location",
      "type": "AWS::S3::ObjectKey",
      "default": "s3://us-east-1.elasticmapreduce.samples/pig-apache-logs/data"
    },
    {
      "id": "myS3OutputLoc",
      "description": "S3 output location",
      "type": "AWS::S3::ObjectKey"
    }
  ]
}

```

```

    }
]
}

```

Note

Sie können diese Objekte der Pipeline-Definitionsdatei anstatt über eine separate Datei auch direkt hinzufügen.

Die folgende Tabelle beschreibt die Attribute für Parameterobjekte.

Parameterattribute

Attribut	Typ	Beschreibung
<code>id</code>	String	Der eindeutige Bezeichner des Parameters. Wenn der Wert bei der Eingabe oder Anzeige maskiert werden soll, fügen Sie als Präfix ein Sternchen ('*') hinzu. Zum Beispiel, <code>*myVariable</code> —. Dabei ist zu beachten, dass der Wert dadurch auch verschlüsselt wird, bevor er durch AWS Data Pipeline gespeichert wird.
<code>description</code>	String	Eine Beschreibung des Parameters.
<code>Typ</code>	Zeichenfolge, Ganzzahl, Double oder AWS::S3::ObjectKey	Der Parametertyp zur Definition des zulässigen Bereichs von Eingabewerten und Validierungsregeln. Der Standardwert ist eine Zeichenfolge.

Attribut	Typ	Beschreibung
optional	Boolesch	Gibt an, ob der Parameter optional oder erforderlich ist. Der Standardwert ist false.
allowedValues	Liste von Zeichenfolgen	Listet alle zulässigen Werte für den Parameter auf.
default	String	Der Standardwert für den Parameter. Wenn Sie mithilfe von Parameterwerten einen Wert für diesen Parameter angeben, wird der Standardwert durch ihn überschrieben.
isArray	Boolesch	Gibt an, ob der Parameter ein Array ist.

Definieren von Parameterwerten

Sie können eine separate Datei zur Definition von Variablen mithilfe von Parameterwerten erstellen. Die folgende JSON-Datei, `file://values.json`, enthält beispielsweise den Wert für eine `myS3OutputLoc` Variable aus der obigen Beispiel-Pipeline-Definition.

```
{
  "values": {
    "myS3OutputLoc": "myOutputLocation"
  }
}
```

Die Pipeline-Definition wird eingereicht

Wenn Sie Ihre Pipeline-Definition senden, können Sie Parameter, Parameterobjekte und Parameterwerte angeben. Sie können den [put-pipeline-definition](#) AWS CLI Befehl beispielsweise wie folgt verwenden:

```
$ aws datapipeline put-pipeline-definition --pipeline-id id --pipeline-definition  
file://pipeline-definition.json \  
--parameter-objects file://parameters.json --parameter-values-uri file://values.json
```

Note

In einer Pipeline-Definition sind maximal 50 Parameter zulässig. Die Größe der Datei für parameter-values-uri darf maximal 15 kB betragen.

Anzeigen Ihrer Pipelines

Sie können Ihre Pipelines über die Befehlszeilenschnittstelle (CLI) anzeigen.

Um Ihre Pipelines anzuzeigen, verwenden Sie AWS CLI

- Verwenden Sie den folgenden [list-pipelines](#)-Befehl, um Ihre Pipelines aufzulisten:

```
aws datapipeline list-pipelines
```

Interpretieren der Pipeline-Statuscodes

Die in der AWS Data Pipeline Konsole und der CLI angezeigten Statusstufen geben den Zustand einer Pipeline und ihrer Komponenten an. Der Pipeline-Status ist vereinfacht ausgedrückt ein Überblick über eine Pipeline. Wenn Sie weitere Informationen benötigen, zeigen Sie den Status der einzelnen Pipeline-Komponenten an.

Der Status einer Pipeline lautet SCHEDULED, wenn sie bereit ist (die Pipeline-Definition hat die Validierung bestanden), aktuell Arbeiten ausführt oder die Ausführung von Arbeiten beendet hat. Der Status einer Pipeline lautet PENDING, wenn sie nicht aktiviert ist oder keine Arbeiten ausführen kann (z. B. hat die Pipeline-Definition die Validierung nicht bestanden).

Eine Pipeline gilt als inaktiv, wenn ihr Status PENDING, INACTIVE oder FINISHED lautet. Für inaktive Pipelines fallen Gebühren an (weitere Informationen finden Sie unter [Preise](#)).

Statuscodes

ACTIVATING

Die Komponente oder Ressource wird gestartet, z. B. eine EC2 Instanz.

CANCELED

Die Komponente wurde von einem Benutzer abgebrochen oder AWS Data Pipeline bevor sie ausgeführt werden konnte. Dies kann automatisch geschehen, wenn ein Fehler in einer anderen Komponente oder Ressource auftritt, von der diese Komponente abhängt.

CASCADE_FAILED

Die Komponente oder Ressource wurde aufgrund eines Kaskadenfehlers aufgrund einer ihrer Abhängigkeiten storniert, aber die Komponente war wahrscheinlich nicht die ursprüngliche Ursache des Fehlers.

DEACTIVATING

Die Pipeline wird deaktiviert.

FAILED

Bei der Komponente oder Ressource ist ein Fehler aufgetreten und sie funktioniert nicht mehr. Wenn eine Komponente oder Ressource ausfällt, kann dies dazu führen, dass Abbrüche und Ausfälle auf andere Komponenten übertragen werden, die von der Komponente oder Ressource abhängen.

FINISHED

Die Komponente hat die ihr zugewiesene Arbeit abgeschlossen.

INACTIVE

Die Pipeline wurde deaktiviert.

PAUSED

Die Komponente wurde angehalten und führt derzeit ihre Arbeit nicht aus.

PENDING

Die Pipeline ist bereit, zum ersten Mal aktiviert zu werden.

RUNNING

Die Ressource läuft und ist bereit, Arbeit zu empfangen.

SCHEDULED

Die Ausführung der Ressource ist geplant.

SHUTTING_DOWN

Die Ressource wird nach erfolgreichem Abschluss ihrer Arbeit heruntergefahren.

SKIPPED

Die Komponente hat Ausführungsintervalle nach der Aktivierung der Pipeline übersprungen, wobei ein Zeitstempel verwendet wurde, der nach dem aktuellen Zeitplan liegt.

TIMEDOUT

Die Ressource hat den `terminateAfter` Schwellenwert überschritten und wurde von AWS Data Pipeline gestoppt. Wenn die Ressource diesen Status erreicht hat, werden die `retryTimeout` Werte `actionOnResourceFailureRetryDelay`, und für diese Ressource AWS Data Pipeline ignoriert. Dieser Status gilt nur für Ressourcen.

VALIDATING

Die Pipeline-Definition wird von validiert AWS Data Pipeline.

WAITING_FOR_RUNNER

Die Komponente wartet darauf, dass ihr Worker-Client ein Arbeitselement abruft. Die Beziehung zwischen Komponente und Worker-Client wird durch die `workerGroup` Felder `runsOn` oder gesteuert, die von dieser Komponente definiert sind.

WAITING_ON_DEPENDENCIES

Die Komponente überprüft, ob ihre standardmäßigen und benutzerdefinierten Voraussetzungen erfüllt sind, bevor sie ihre Arbeit ausführt.

Interpretieren des Pipeline- und Komponenten-Zustands

Jede Pipeline und Komponente innerhalb der betreffenden Pipeline geben als Zustand `HEALTHY`, `ERROR`, `"-"`, `No Completed Executions` oder `No Health Information Available` zurück. Für eine Pipeline wird erst dann ein Zustand angegeben, nachdem eine Pipeline-Komponente erstmals ausgeführt wurde oder wenn die Vorbedingungen einer Komponente fehlgeschlagen sind. Der Zustand einzelner Komponenten wird im Pipeline-Zustand aggregiert. Der Fehlerstatus ist erst sichtbar, wenn Sie Ihre Pipeline-Ausführungsdetails anzeigen.

Pipeline-Zustände

HEALTHY

Der aggregierte Zustand aller Komponenten ist **HEALTHY**. Dies bedeutet, dass mindestens eine Komponente erfolgreich abgeschlossen worden sein muss. Sie können auf den Status **HEALTHY** klicken, um die neueste erfolgreich abgeschlossene Pipeline-Komponenten-Instance auf der Seite **Execution Details** zu sehen.

ERROR

Mindestens eine Komponente in der Pipeline verfügt über den Zustand **ERROR**. Sie können auf den Status **ERROR** klicken, um die neueste fehlgeschlagene Pipeline-Komponenten-Instance auf der Seite **Execution Details** zu sehen.

No Completed Executions oder No Health Information Available.

Für diese Pipeline wurde kein Zustand gemeldet.

Note

Während Komponenten ihren Zustand beinahe sofort aktualisieren, kann es bis zu fünf Minuten dauern, bis der Pipeline-Zustand aktualisiert wird.

Komponenten-Zustand

HEALTHY

Der Zustand einer Komponente (Activity oder DataNode) ist **HEALTHY**, wenn sie nach einer erfolgreichen Ausführung mit dem Status **FINISHED** oder **MARK_FINISHED** markiert wurde. Sie können auf den Namen der Komponente oder auf den Status **HEALTHY** klicken, um auf der Seite **Execution Details** die neuesten erfolgreich abgeschlossenen Pipeline-Komponenten-Instances anzuzeigen.

ERROR

Es ist ein Fehler auf Komponentenebene aufgetreten oder eine ihrer Vorbedingungen ist fehlgeschlagen. Dieser Fehler wird durch einen Status von **FAILED**, **TIMEOUT** oder **CANCELED** ausgelöst. Sie können auf den Namen der Komponente oder auf den Status **ERROR** klicken, um auf der Seite **Execution Details** die neueste fehlgeschlagene Pipeline-Komponenten-Instance anzuzeigen.

No Completed Executions oder No Health Information Available

Für diese Komponente wurde kein Zustand gemeldet.

Anzeigen Ihrer Pipeline-Definitionen

Verwenden Sie die Befehlszeilenschnittstelle (CLI), um Ihre Pipeline-Definition anzuzeigen. Die CLI druckt eine Pipeline-Definitionsdatei im JSON-Format. Weitere Informationen zur Syntax und Nutzung von Pipeline-Definitionsdateien finden Sie unter [Syntax der Pipeline-Definitionsdatei](#).

Wenn Sie die CLI verwenden, empfiehlt es sich, die Pipeline-Definition abzurufen, bevor Sie Änderungen einreichen, da es möglich ist, dass ein anderer Benutzer oder Prozess die Pipeline-Definition geändert hat, nachdem Sie zuletzt damit gearbeitet haben. Durch Herunterladen einer Kopie der aktuellen Definition, auf der Sie Ihre Änderungen basieren, können Sie sicher sein, dass Sie an der neuesten Pipeline-Definition arbeiten. Auch wird empfohlen, die Pipeline-Definition nach den vorgenommenen Änderungen erneut abzurufen, um sicherzugehen, dass die Aktualisierung erfolgreich war.

Wenn Sie die CLI verwenden, können Sie zwei verschiedene Versionen Ihrer Pipeline erhalten. Die Version `active` ist die Pipeline, die gerade ausgeführt wird. Die Version `latest` ist eine Kopie, die beim Bearbeiten einer gerade ausgeführten Pipeline erstellt wird. Wenn Sie die bearbeitete Pipeline hochladen, ist sie nun die Version `active` und die vorherige Version `active` ist nicht mehr verfügbar.

Um eine Pipeline-Definition mit dem zu erhalten AWS CLI

Um die vollständige Pipeline-Definition abzurufen, verwenden Sie den [get-pipeline-definition](#)Befehl. Die Pipeline-Definition wird zur Standardausgabe (stdout) ausgegeben.

Das folgende Beispiel ruft die Pipeline-Definition für die angegebene Pipeline ab.

```
aws datapipeline get-pipeline-definition --pipeline-id df-00627471S0VYZEXAMPLE
```

Zum Abrufen einer bestimmten Version einer Pipeline verwenden Sie die Option `--version`. Das folgende Beispiel ruft die Version `active` der angegebenen Pipeline ab.

```
aws datapipeline get-pipeline-definition --version active --id df-00627471S0VYZEXAMPLE
```

Anzeigen von Pipeline-Instance Details

Sie können den Fortschritt Ihrer Pipeline überwachen. Weitere Informationen zum Instance-Status finden Sie unter [Interpretieren der Pipeline-Statusdetails](#). Weitere Informationen zur Fehlerbehebung bei fehlgeschlagenen oder unvollständigen Instance-Ausführungen Ihrer Pipeline finden Sie unter [Beheben typischer Probleme](#).

Um den Fortschritt einer Pipeline mit dem zu überwachen AWS CLI

Um Pipeline-Instance-Details abzurufen, wie z. B. wie oft eine Pipeline ausgeführt wurde, verwenden Sie den Befehl [list-runs](#). Mit diesem Befehl können Sie die Liste der zurückgegebenen Ausführungen nach ihrem aktuellem Status oder dem Datumsbereich filtern, in dem sie gestartet wurden. Das Filtern der Ergebnisse ist hilfreich, da der Ausführungsverlauf je nach Alter und Zeitplanung der Pipeline groß sein kann.

Das folgende Beispiel ruft Informationen zu allen Ausführungen ab.

```
aws datapipeline list-runs --pipeline-id df-00627471S0VYZEXAMPLE
```

Das folgende Beispiel ruft Informationen zu allen abgeschlossenen Ausführungen ab.

```
aws datapipeline list-runs --pipeline-id df-00627471S0VYZEXAMPLE --status finished
```

Das folgende Beispiel ruft Informationen zu allen in einem bestimmten Zeitrahmen gestarteten Ausführungen ab.

```
aws datapipeline list-runs --pipeline-id df-00627471S0VYZEXAMPLE --start-interval  
"2013-09-02", "2013-09-11"
```

Anzeigen von Pipeline-Protokollen

Die Protokollierung auf Pipeline-Ebene wird bei der Pipeline-Erstellung unterstützt, indem ein Amazon S3 S3-Speicherort entweder in der Konsole oder mit einem `pipelineLogUri` im Standardobjekt in SDK/CLI angegeben wird. Die Verzeichnisstruktur für jede Pipeline in diesem URL wird nachstehend beschrieben:

```
pipelineId  
  -componentName
```

```
-instanceId
-attemptId
```

Für die Pipeline df-00123456ABC7DEF8HIJK sieht die Verzeichnisstruktur folgendermaßen aus:

```
df-00123456ABC7DEF8HIJK
  -ActivityId_fxNzc
    -@ActivityId_fxNzc_2014-05-01T00:00:00
      -@ActivityId_fxNzc_2014-05-01T00:00:00_Attempt=1
```

Für ShellCommandActivity werden Protokolle für stderr und stdout, die diesen Aktivitäten zugeordnet sind, bei jedem Versuch im Verzeichnis gespeichert.

Für Ressourcen wie EmrCluster mit festgelegtem emrLogUri hat dieser Wert Vorrang. Andernfalls folgen die Ressourcen (einschließlich der TaskRunner Protokolle für diese Ressourcen) der obigen Struktur der Pipeline-Protokollierung.

So zeigen Sie Protokolle für einen bestimmten Pipeline-Lauf an:

1. Rufen Sie das ab, ObjectId indem query-objects Sie die genaue Objekt-ID aufrufen. Zum Beispiel:

```
aws datapipeline query-objects --pipeline-id <pipeline-id> --sphere ATTEMPT --region ap-northeast-1
```

query-objects ist eine paginierte CLI und kann ein Paginierungstoken zurückgeben, wenn es mehr Ausführungen für die angegebene gibt. pipeline-id Sie können das Token verwenden, um alle Versuche durchzuführen, bis Sie das erwartete Objekt gefunden haben. Ein zurückgegebenes Objekt ObjectId würde beispielsweise wie folgt aussehen:@TableBackupActivity_2023-05-020T18:05:18_Attempt=1.

2. Rufen Sie mit dem ObjectId den Speicherort des Protokolls ab mit:

```
aws datapipeline describe-objects --pipeline-id <pipeline-id> --object-ids <object-id>
  --query "pipelineObjects[].fields[?key=='@logLocation'].stringValue"
```

Fehlermeldung einer fehlgeschlagenen Aktivität

Um die Fehlermeldung zu erhalten, rufen Sie zuerst die ObjectId Verwendung abquery-objects.

Verwenden Sie nach dem Abrufen der ObjectID fehlgeschlagenen Datei die `describe-objects` CLI, um die eigentliche Fehlermeldung abzurufen.

```
aws datapipeline describe-objects --region ap-northeast-1 --pipeline-id <pipeline-id> --object-ids <object-id> --query "pipelineObjects[].fields[?key=='errorMessage'].stringValue"
```

Brechen Sie ein Objekt ab, führen Sie es erneut aus oder markieren Sie es als abgeschlossen

Verwenden Sie die `set-status` CLI, um ein laufendes Objekt abzubrechen, ein fehlgeschlagenes Objekt erneut auszuführen oder ein laufendes Objekt als Fertig zu markieren.

Rufen Sie zunächst die Objekt-ID mit der `query-objects` CLI ab. Zum Beispiel:

```
aws datapipeline query-objects --pipeline-id <pipeline-id> --sphere INSTANCE --region ap-northeast-1
```

Verwenden Sie die `set-status` CLI, um den Status des gewünschten Objekts zu ändern. Zum Beispiel:

```
aws datapipeline set-status --pipeline-id <pipeline-id> --region ap-northeast-1 --status TRY_CANCEL --object-ids <object-id>
```

Bearbeiten Ihrer Pipeline

Wenn Sie bestimmte Aspekte einer Ihrer Pipelines ändern müssen, können Sie die Definition der Pipeline entsprechend aktualisieren. Wenn Änderungen an einer derzeit ausgeführten Pipeline vorgenommen wurden, müssen Sie die Pipeline erneut aktivieren, damit die Änderungen wirksam werden. Zudem ist es möglich, eine oder mehrere Pipeline-Komponenten erneut auszuführen.

Inhalt

- [Einschränkungen](#)
- [Bearbeiten einer Pipeline mit dem AWS CLI](#)

Einschränkungen

Solange sich die Pipeline im PENDING Status befindet und nicht aktiviert ist, können Sie keine Änderungen daran vornehmen. Nachdem eine Pipeline aktiviert wurde, gelten beim Bearbeiten der

Pipeline die folgenden Einschränkungen. Die von Ihnen vorgenommenen Änderungen werden für neue Ausführungen der Pipeline-Objekte übernommen, nachdem Sie sie speichern und die Pipeline dann erneut aktivieren.

- Ein Objekt kann nicht entfernt werden
- Der Planungszeitraum eines vorhandenen Objekts kann nicht geändert werden
- Referenzfelder in einem vorhandenen Objekt können nicht hinzugefügt, gelöscht oder abgeändert werden
- In dem Ausgabefeld eines neuen Objekts kann nicht auf ein vorhandenes Objekt verwiesen werden
- Das geplante Anfangsdatum eines Objekts kann nicht geändert werden (aktivieren Sie stattdessen die Pipeline mit bestimmten Angaben für Datum und Uhrzeit)

Bearbeiten einer Pipeline mit dem AWS CLI

Sie können eine Pipeline mit Befehlszeilen-Tools bearbeiten.

Laden Sie zunächst mit dem [get-pipeline-definition](#)Befehl eine Kopie der aktuellen Pipeline-Definition herunter. Auf diese Weise können Sie sicher sein, dass Sie die neuesten Pipeline-Definition abändern. In dem folgenden Beispiel wird die Pipeline-Definition zur Standardausgabe (stdout) ausgegeben.

```
aws datapipeline get-pipeline-definition --pipeline-id df-00627471S0VYZEXAMPLE
```

Speichern Sie die Pipeline-Definition in einer Datei und bearbeiten Sie sie nach Bedarf. Aktualisieren Sie Ihre Pipeline-Definition mit dem [put-pipeline-definition](#)Befehl. Im folgenden Beispiel wird die Pipeline-Definitionsdatei hochgeladen.

```
aws datapipeline put-pipeline-definition --pipeline-id df-00627471S0VYZEXAMPLE --  
pipeline-definition file://MyEmrPipelineDefinition.json
```

Sie können die Pipeline-Definition mit dem Befehl `get-pipeline-definition` erneut abrufen, um sicherzustellen, dass die Aktualisierung erfolgreich war. Verwenden Sie zum Aktivieren der Pipeline den folgenden [activate-pipeline](#)-Befehl:

```
aws datapipeline activate-pipeline --pipeline-id df-00627471S0VYZEXAMPLE
```

Wenn Sie möchten, können Sie die Pipeline ab einem bestimmten Datum und einer bestimmten Uhrzeit aktivieren. Verwenden Sie hierzu folgendermaßen die Option `--start-timestamp`:

```
aws datapipeline activate-pipeline --pipeline-id df-00627471S0VYZEXAMPLE --start-timestamp YYYY-MM-DDTHH:MM:SSZ
```

Mit dem Befehl [set-status](#) können Sie eine oder mehrere Pipeline-Komponenten erneut ausführen.

Klonen Ihrer Pipeline

Klonen erstellt eine Kopie einer Pipeline und ermöglicht Ihnen einen Namen für die neue Pipeline anzugeben. Pipelines können in jedem Zustand geklont werden, sogar mit Fehlern. Die neue Pipeline verbleibt jedoch im Zustand PENDING, bis Sie sie manuell aktivieren. Der Clone-Vorgang verwendet für die neue Pipeline die neueste Version der ursprünglichen Pipeline-Definition und nicht die aktive Version. Bei dem Klonvorgang wird nicht der vollständige Zeitplan aus der ursprünglichen Pipeline in die neue Pipeline kopiert, sondern nur der festgelegte Zeitraum.

So klonen Sie eine Pipeline mit der AWS CLI:

1. Erstellen Sie eine neue Pipeline mit einem neuen Namen und einer eindeutigen ID. Notieren Sie sich die zurückgegebene Pipeline-ID.
2. Verwenden Sie die `get-pipeline-definition` CLI, um die Pipeline-Definition der vorhandenen Pipeline abzurufen, die geklont werden soll, und sie in eine temporäre Datei zu schreiben. Notieren Sie sich den absoluten Pfad der Datei.
3. Verwenden Sie die `put-pipeline-definition` CLI, um die Pipeline-Definition von der vorhandenen Pipeline in die neue Pipeline zu kopieren.
4. Verwenden Sie die `get-pipeline-definition` CLI, um die Definition der neuen Pipeline abzurufen und die Pipeline-Definition zu überprüfen.

```
# Create Pipeline (returns <new-pipeline-id>)
aws datapipeline create-pipeline --name my-cloned-pipeline --unique-id my-cloned-pipeline --region ap-northeast-1

#Get pipeline definition of existing pipeline
aws datapipeline get-pipeline-definition --pipeline-id <existing-pipeline-id> --region ap-northeast-1 > existing_pipeline_definition.json

# Put pipeline definition to new pipeline
```

```
aws datapipeline put-pipeline-definition --pipeline-id <new-pipeline-id> --region ap-northeast-1 --pipeline-definition file://<absolute_path_to_existing_pipeline_definition.json>

# get pipeline definition of new pipeline
aws datapipeline get-pipeline-definition --pipeline-id <new-pipeline-id> --region ap-northeast-1
```

Tagging Ihrer Pipeline

Tags sind Schlüssel/Wert-Paare mit Unterscheidung nach Groß-/Kleinschreibung, die aus einem Schlüssel und einem optionalen Wert bestehen, die beide vom Benutzer definiert wurden. Sie können auf jede Pipeline bis zu 10 Tags anwenden. Tag-Schlüssel müssen für jede Pipeline eindeutig sein. Wenn Sie ein Tag mit einem Schlüssel hinzufügen, der der Pipeline bereits zugeordnet ist, ändert sich der Wert dieses Tags.

Durch das Anwenden eines Tags auf eine Pipeline werden die Tags auch an die zugrunde liegenden Ressourcen weitergegeben (z. B. Amazon EMR-Cluster und EC2 Amazon-Instances). Diese Tags werden allerdings nicht auf Ressourcen im Zustand FINISHED oder in einem anderen beendeten Zustand übertragen. Sie können mit der CLI bei Bedarf Tags auf diese Ressourcen anwenden.

Wenn Sie ein Tag nicht mehr benötigen, können Sie es von Ihrer Pipeline entfernen.

So versehen Sie Ihre Pipeline mit Tags über die AWS CLI

Um einer neuen Pipeline Tags hinzuzufügen, fügen Sie die `--tags` Option zu Ihrem Befehl [create-pipeline](#) hinzu. Die folgende Option erstellt z. B. eine Pipeline mit zwei Tags, einem `environment`-Tag mit dem Wert `production` und einem `owner`-Tag mit dem Wert `sales`.

```
--tags key=environment,value=production key=owner,value=sales
```

Um Tags zu einer vorhandenen Pipeline hinzuzufügen, verwenden Sie den Befehl [add-tags](#) folgendermaßen:

```
aws datapipeline add-tags --pipeline-id df-00627471SOVYZEXAMPLE --tags key=environment,value=production key=owner,value=sales
```

Um Tags von einer vorhandenen Pipeline zu entfernen, verwenden Sie den Befehl [remove-tags](#) folgendermaßen:

```
aws datapipeline remove-tags --pipeline-id df-00627471SOVYZEXAMPLE --tag-keys environment owner
```

Deaktivieren Ihrer Pipeline

Wenn eine gerade ausgeführte Pipeline deaktiviert wird, wird die Ausführung der Pipeline angehalten. Um die Ausführung der Pipeline fortzusetzen, können Sie die Pipeline aktivieren. Dies ermöglicht Ihnen, die folgenden Änderungen vorzunehmen. Wenn Sie beispielsweise Daten in eine Datenbank schreiben, für die eine Wartung geplant ist, können Sie die Pipeline deaktivieren, und warten, bis die Wartung abgeschlossen ist, und die Pipeline dann aktivieren.

Beim Deaktivieren der Pipeline können Sie den Umgang mit laufenden Aktivitäten festlegen. Standardmäßig werden diese Aktivitäten sofort abgebrochen. Alternativ können Sie AWS Data Pipeline veranlassen, mit dem Deaktivieren der Pipeline zu warten, bis die Aktivitäten abgeschlossen wurden.

Beim Aktivieren einer deaktivierten Pipeline können Sie festlegen, wann sie fortgesetzt wird. Wenn Sie die API AWS CLI oder die API verwenden, wird die Pipeline standardmäßig mit der letzten abgeschlossenen Ausführung fortgesetzt, oder Sie können das Datum und die Uhrzeit angeben, zu der die Pipeline wieder aufgenommen werden soll.

Deaktivieren Ihrer Pipeline über die AWS CLI

Verwenden Sie den folgenden [deactivate-pipeline](#)-Befehl zum Deaktivieren einer Pipeline:

```
aws datapipeline deactivate-pipeline --pipeline-id df-00627471SOVYZEXAMPLE
```

Um die Pipeline erst nach der Ausführung aller Aktivitäten zu deaktivieren, fügen Sie die Option `--no-cancel-active` wie folgt hinzu:

```
aws datapipeline deactivate-pipeline --pipeline-id df-00627471SOVYZEXAMPLE --no-cancel-active
```

Wenn Sie bereit sind, können Sie die Ausführung der Pipeline mit dem folgenden [activate-pipeline](#)-Befehl an der Stelle fortsetzen, an der sie unterbrochen wurde:

```
aws datapipeline activate-pipeline --pipeline-id df-00627471SOVYZEXAMPLE
```

Um die Pipeline ab einem bestimmten Datum und einer bestimmten Uhrzeit zu starten, fügen Sie die Option `--start-timestamp` wie folgt hinzu:

```
aws datapipeline activate-pipeline --pipeline-id df-00627471S0VYZEXAMPLE --start-timestamp YYYY-MM-DDTHH:MM:SSZ
```

Löschen Ihrer Pipeline

Wenn Sie eine Pipeline nicht mehr benötigen, z. B. eine Pipeline, die beim Testen der Anwendung erstellt wurde, sollten Sie sie löschen, damit sie nicht mehr aktiv genutzt wird. Beim Löschen wird eine Pipeline in den gelöschten Zustand versetzt. Im gelöschten Zustand besitzen Pipelines keine Pipeline-Definition und keinen Ausführungsverlauf mehr. Für solche Pipeline können daher keine Vorgänge mehr durchgeführt werden, einschließlich deren Beschreibung.

 **Important**

Gelöschte Pipelines können nicht wiederhergestellt werden. Stellen Sie daher sicher, dass Sie die Pipeline zukünftig nicht mehr benötigen, bevor Sie sie löschen.

Um eine Pipeline mit dem zu löschen AWS CLI

Verwenden Sie zum Löschen einer Pipeline den Befehl [delete-pipeline](#). Der folgende Befehl löscht die angegebene Pipeline.

```
aws datapipeline delete-pipeline --pipeline-id df-00627471S0VYZEXAMPLE
```

Staging von Daten und Tabellen mit Pipeline-Aktivitäten

AWS Data Pipeline kann Eingabe- und Ausgabedaten in Ihren Pipelines bereitstellen, um die Verwendung bestimmter Aktivitäten wie `ShellCommandActivity` und `HiveActivity` zu vereinfachen.

Daten-Staging gibt Ihnen die Möglichkeit, Daten von einem Eingabedatenknoten in die Ressource zu kopieren, die die Aktivität ausführt, und anschließend genauso von der Resource in den Ausgabedatenknoten.

Die bereitgestellten Daten auf der Amazon EMR- oder EC2 Amazon-Ressource sind mithilfe spezieller Variablen in den Shell-Befehlen oder Hive-Skripten der Aktivität verfügbar.

Tabellen-Staging ist mit Daten-Staging vergleichbar, außer dass die Daten speziell in Form von Datenbanktabellen bereitgestellt werden.

AWS Data Pipeline unterstützt die folgenden Staging-Szenarien:

- Daten-Staging mit `ShellCommandActivity`
- Tabellen-Staging mit Hive und zum Staging fähigen Datenknoten
- Tabellen-Staging mit Hive und nicht zum Staging fähigen Datenknoten

 Note

Staging funktioniert nur, wenn das Feld `stage` für eine Aktivität wie etwa `ShellCommandActivity` auf `true` festgelegt ist. Weitere Informationen finden Sie unter [ShellCommandActivity](#).

Darüber hinaus sind zwischen Datenknoten und Aktivitäten vier verschiedene Beziehungen möglich:

Staging von Daten lokal auf einer Ressource

Die Eingabedaten werden automatisch in das lokale Dateisystem der Ressource kopiert. Die Ausgabedaten werden automatisch aus dem lokalen Dateisystem auf den Ausgabedatenknoten kopiert. Wenn Sie z. B. `ShellCommandActivity`-Eingaben und -Ausgaben mit `Staging = true` konfigurieren, sind die Eingabedaten als `INPUTx_STAGING_DIR` und die Ausgabedaten als `OUTPUTx_STAGING_DIR` verfügbar, wobei `x` die Anzahl der Ein- oder Ausgaben ist.

Staging von Eingabe- und -Ausgabedefinitionen für eine Aktivität

Das Format der Eingabedaten (Spaltennamen und Tabellennamen) wird automatisch auf die Ressource der Aktivität kopiert. Beispielsweise, wenn Sie `HiveActivity` mit `Staging = true` konfigurieren. Das auf dem Eingabedatenknoten `S3DataNode` angegebene Datenformat wird zum Bereitstellen der Tabellendefinition aus der Hive-Tabelle verwendet.

Staging nicht aktiviert

Die Eingabe- und Ausgabe-Objekte und ihre Felder sind für die Aktivität zwar verfügbar, die Daten selbst jedoch nicht. Dies trifft z. B. bei `EmrActivity` standardmäßig zu. Andere Aktivitäten müssen Sie dazu mit `Staging = false` konfigurieren. In dieser Konfiguration stehen die Datenfelder der Aktivität zur Verfügung, um mithilfe der AWS Data Pipeline Ausdruckssyntax auf sie zu verweisen. Dies geschieht nur, wenn die Abhängigkeit erfüllt ist. Dies dient alleinig

der Abhängigkeitsprüfung. Der Code in der Aktivität veranlasst das Kopieren der Daten von der Eingabe auf die Ressource, auf der die Aktivität ausgeführt wird.

Abhängigkeitsbeziehung zwischen Objekten

Zwischen zwei Objekten besteht eine Abhängigkeitsbeziehung, die zu einer ähnlichen Situation führt, wie wenn Staging nicht aktiviert ist. In diesem Fall fungiert ein Datenknoten oder eine Aktivität als Vorbedingung für die Ausführung einer anderen Aktivität.

Daten-Staging mit ShellCommandActivity

Stellen Sie sich ein Szenario vor, in ShellCommandActivity dem S3DataNode Objekte mit einem Objekt als Dateneingabe und -ausgabe verwendet werden. AWS Data Pipeline stellt die Datenknoten automatisch bereit, um sie mithilfe der Umgebungsvariablen für den Shell-Befehl zugänglich zu machen, als ob es sich um lokale Dateiordner \${INPUT1_STAGING_DIR} handeln würde, \${OUTPUT1_STAGING_DIR} wie im folgenden Beispiel gezeigt. Der numerische Teil der Variablen namens INPUT1_STAGING_DIR und OUTPUT1_STAGING_DIR wird abhängig von der Anzahl der Datenknoten erhöht, auf die Ihre Aktivität verweist.

Note

Dieses Szenario funktioniert nur wie beschrieben, wenn Ihre Dateneingabe- und -ausgaben S3DataNode-Objekte sind. Außerdem ist das Staging von Ausgabedaten nur zulässig, wenn für das S3DataNode-Ausgabeobjekt directoryPath festgelegt ist.

```
{  
  "id": "AggregateFiles",  
  "type": "ShellCommandActivity",  
  "stage": "true",  
  "command": "cat ${INPUT1_STAGING_DIR}/part* > ${OUTPUT1_STAGING_DIR}/aggregated.csv",  
  "input": {  
    "ref": "MyInputData"  
  },  
  "output": {  
    "ref": "MyOutputData"  
  }  
},  
{  
  "id": "MyInputData",  
  "type": "S3DataNode",  
  "stage": "true",  
  "path": "s3://my-bucket/input/part*",  
  "region": "us-west-2",  
  "format": "CSV",  
  "header": true,  
  "compression": "none",  
  "partitionKeys": []  
},  
{  
  "id": "MyOutputData",  
  "type": "S3DataNode",  
  "stage": "true",  
  "path": "s3://my-bucket/output/aggregated.csv",  
  "region": "us-west-2",  
  "format": "CSV",  
  "header": false,  
  "compression": "none",  
  "partitionKeys": []  
}
```

```
"type": "S3DataNode",
"schedule": {
    "ref": "MySchedule"
},
"filePath": "s3://my_bucket/source/#{format(@scheduledStartTime, 'YYYY-MM-
dd_HH:mm:ss')}/items"
},
{
    "id": "MyOutputData",
    "type": "S3DataNode",
    "schedule": {
        "ref": "MySchedule"
    },
    "directoryPath": "s3://my_bucket/destination/#{format(@scheduledStartTime, 'YYYY-MM-
dd_HH:mm:ss')}"
}
},
...
...
```

Tabellen-Staging mit Hive und zum Staging fähigen Datenknoten

Stellen Sie sich ein Szenario vor, in `HiveActivity` dem `S3DataNode` Objekte mit einem Objekt als Dateneingabe und -ausgabe verwendet werden. AWS Data Pipeline stellt die Datenknoten automatisch bereit, um sie für das Hive-Skript zugänglich zu machen, als ob es sich um Hive-Tabellen handeln würde, die die Variablen verwenden `#{input1}` und `#{output1}` wie im folgenden Beispiel für gezeigt. `HiveActivity` Der numerische Teil der Variablen namens `input` und `output` wird abhängig von der Anzahl der Datenknoten erhöht, auf die Ihre Aktivität verweist.

Note

Dieses Szenario funktioniert nur wie beschrieben, wenn die Dateneingaben und -ausgaben `S3DataNode`- oder `MySQLDataNode`-Objekte sind. Tabelle-Staging wird für `DynamoDBDataNode` nicht unterstützt.

```
{
    "id": "MyHiveActivity",
    "type": "HiveActivity",
    "schedule": {
        "ref": "MySchedule"
    }
}
```

```
},
"runsOn": {
    "ref": "MyEmrResource"
},
"input": {
    "ref": "MyInputData"
},
"output": {
    "ref": "MyOutputData"
},
"hiveScript": "INSERT OVERWRITE TABLE ${output1} select * from ${input1};"
},
{
    "id": "MyInputData",
    "type": "S3DataNode",
    "schedule": {
        "ref": "MySchedule"
    },
    "directoryPath": "s3://test-hive/input"
}
},
{
    "id": "MyOutputData",
    "type": "S3DataNode",
    "schedule": {
        "ref": "MySchedule"
    },
    "directoryPath": "s3://test-hive/output"
}
},
...
...
```

Tabellen-Staging mit Hive und nicht zum Staging fähigen Datenknoten

Angenommen, eine `HiveActivity` wird für die Dateneingabe mit `DynamoDBDataNode` und für die Datenausgabe mit einem `S3DataNode`-Objekt verwendet. Für ist kein Daten-Staging verfügbar `DynamoDBDataNode`, daher müssen Sie die Tabelle zunächst manuell in Ihrem Hive-Skript erstellen und dabei den Variablennamen verwenden, um auf die `DynamoDB`-Tabelle `#{input.tableName}` zu verweisen. Eine ähnliche Nomenklatur gilt, wenn die `DynamoDB`-Tabelle die Ausgabe ist, es sei denn, Sie verwenden eine Variable. `#{output.tableName}` Da in diesem Beispiel für das Ausgabeobjekt `S3DataNode` Staging verfügbar ist, können Sie unter `#{output1}` auf den Ausgabedatenknoten verweisen.

Note

In diesem Beispiel hat die Variable für den Tabellennamen das Präfix # (Hash), da Ausdrücke für den Zugriff auf oder AWS Data Pipeline verwendet werden. `tableName` `directoryPath` Weitere Hinweise zur Funktionsweise der Ausdrucksauswertung in AWS Data Pipeline finden Sie unter [Ausdrucksauswertung](#).

```
{  
    "id": "MyHiveActivity",  
    "type": "HiveActivity",  
    "schedule": {  
        "ref": "MySchedule"  
    },  
    "runsOn": {  
        "ref": "MyEmrResource"  
    },  
    "input": {  
        "ref": "MyDynamoData"  
    },  
    "output": {  
        "ref": "MyS3Data"  
    },  
    "hiveScript": "-- Map DynamoDB Table  
SET dynamodb.endpoint=dynamodb.us-east-1.amazonaws.com;  
SET dynamodb.throughput.read.percent = 0.5;  
CREATE EXTERNAL TABLE dynamodb_table (item map<string,string>)  
STORED BY 'org.apache.hadoop.hive.dynamodb.DynamoDBStorageHandler'  
TBLPROPERTIES ("dynamodb.table.name" = "#{input.tableName}");  
INSERT OVERWRITE TABLE ${output1} SELECT * FROM dynamodb_table;"  
},  
{  
    "id": "MyDynamoData",  
    "type": "DynamoDBDataNode",  
    "schedule": {  
        "ref": "MySchedule"  
    },  
    "tableName": "MyDDBTable"  
},  
{  
    "id": "MyS3Data",  
    "type": "S3DataNode",  
}
```

```
"schedule": {  
    "ref": "MySchedule"  
},  
"directoryPath": "s3://test-hive/output"  
}  
},  
...  
}
```

Verwenden einer Pipeline mit Ressourcen in mehreren Regionen

Standardmäßig werden die `EmrCluster` Ressourcen `Ec2Resource` und in derselben Region ausgeführt, AWS Data Pipeline unterstützen jedoch die Fähigkeit AWS Data Pipeline, Datenflüsse über mehrere Regionen hinweg zu orchestrieren, z. B. die Ausführung von Ressourcen in einer Region, die Eingabedaten aus einer anderen Region konsolidieren. Da Sie Ressourcen die Ausführung in einer bestimmten Region erlauben können, haben Sie auch die Flexibilität, Ressourcen zusammen mit ihren abhängigen Datensätzen in derselben Region anzusiedeln und die Leistung zu maximieren, indem Latenzen verringert und Kosten für regionsübergreifende Datenübertragungen vermieden werden. Sie können Ressourcen so konfigurieren, dass sie in einer anderen Region ausgeführt AWS Data Pipeline werden, als wenn Sie das `region` Feld für `Ec2Resource` und `EmrCluster` verwenden.

Die folgende Beispiel-Pipeline-JSON-Datei zeigt, wie eine `EmrCluster` Ressource in der Region Europa (Irland) ausgeführt wird, wobei davon ausgegangen wird, dass in derselben Region eine große Datenmenge vorhanden ist, an der der Cluster arbeiten kann. In diesem Beispiel ist der einzige Unterschied zu einer typischen Pipeline, dass der Wert des Feldes `region` für `EmrCluster` auf `eu-west-1` eingestellt ist.

```
{  
    "objects": [  
        {  
            "id": "Hourly",  
            "type": "Schedule",  
            "startDateTime": "2014-11-19T07:48:00",  
            "endDateTime": "2014-11-21T07:48:00",  
            "period": "1 hours"  
        },  
        {  
            "id": "MyCluster",  
            "type": "EmrCluster",  
            "masterInstanceType": "m3.medium",  
            "region": "eu-west-1"  
        }  
    ]  
}
```

```
"region": "eu-west-1",
"schedule": {
    "ref": "Hourly"
},
{
    "id": "MyEmrActivity",
    "type": "EmrActivity",
    "schedule": {
        "ref": "Hourly"
    },
    "runsOn": {
        "ref": "MyCluster"
    },
    "step": "/home/hadoop/contrib/streaming/hadoop-streaming.jar,-input,s3n://elasticmapreduce/samples/wordcount/input,-output,s3://eu-west-1-bucket/wordcount/output/#{@scheduledStartTime},-mapper,s3n://elasticmapreduce/samples/wordcount/wordSplitter.py,-reducer,aggregate"
}
]
}
```

In der folgenden Tabelle werden die Regionen aufgelistet, die Sie auswählen können, sowie die im Feld `region` zu verwendenden zugehörigen Regionscodes.

Note

Die folgende Liste enthält Regionen, in denen Workflows orchestriert und Amazon EMR- oder EC2 Amazon-Ressourcen gestartet werden AWS Data Pipeline können. AWS Data Pipeline wird in diesen Regionen möglicherweise nicht unterstützt. Informationen zu Regionen, in denen dies unterstützt AWS Data Pipeline wird, finden Sie unter [AWS-Regionen und Endpunkte](#).

Name der Region	Regionscode
USA Ost (Nord-Virginia)	us-east-1
USA Ost (Ohio)	us-east-2
USA West (Nordkalifornien)	us-west-1

Name der Region	Regionscode
USA West (Oregon)	us-west-2
Kanada (Zentral)	ca-central-1
Europa (Irland)	eu-west-1
Europa (London)	eu-west-2
Europa (Frankfurt)	eu-central-1
Asien-Pazifik (Singapur)	ap-southeast-1
Asien-Pazifik (Sydney)	ap-southeast-2
Asien-Pazifik (Mumbai)	ap-south-1
Asien-Pazifik (Tokio)	ap-northeast-1
Asien-Pazifik (Seoul)	ap-northeast-2
Südamerika (São Paulo)	sa-east-1

Cascading-Ausfälle und erneute Ausführungen

AWS Data Pipeline ermöglicht es Ihnen, das Verhalten von Pipeline-Objekten zu konfigurieren, wenn eine Abhängigkeit fehlschlägt oder von einem Benutzer gestrichen wird. Sie können sicherstellen, dass Ausfälle zu anderen Pipeline-Objekten (Verbrauchern) kaskadieren, um ein unendlich langes Warten zu verhindern. Für alle Aktivitäten, Datenknoten und Vorbedingungen ist ein Feld namens `failureAndRerunMode` mit dem Standardwert `none` vorhanden. Um das Kaskadieren von Ausfällen zu aktivieren, stellen Sie das Feld `failureAndRerunMode` auf `cascade` ein.

Wenn dieses Feld aktiviert ist, treten kaskadierende Ausfälle auf, wenn ein Pipeline-Objekt im Zustand `WAITING_ON_DEPENDENCIES` blockiert wird und alle Abhängigkeiten ohne ausstehenden Befehl fehlgeschlagen sind. Während eines kaskadierenden Ausfalls treten die folgenden Ereignisse ein:

- Wenn ein Objekt fehlschlägt, wird für seine Verbraucher `CASCADE_FAILED` eingestellt und für die Vorbedingungen des Originalobjekts und seiner Verbraucher wird `CANCELED` festgelegt.

- Alle Objekte, die sich bereits im Zustand FINISHED, FAILED oder CANCELED befinden, werden ignoriert.

Das Kaskadieren von Ausfällen funktioniert nicht bei Abhängigkeiten (stromaufwärts) eines ausgefallenen Objekts, außer bei zugehörigen Vorbedingungen des ursprünglichen ausgefallenen Objekts. Pipeline-Objekte, die vom Kaskadieren eines Ausfalls betroffen sind, können beliebig viele Wiederholungen oder Post-Aktionen, wie z. B. onFail, auslösen.

Die detaillierten Auswirkungen eines kaskadierenden Ausfalls sind vom Objekttyp abhängig.

Aktivitäten

Der Zustand einer Aktivität ändert sich in CASCADE_FAILED, wenn irgendwelche seiner Abhängigkeiten ausfallen, und es löst anschließend einen kaskadierenden Ausfall bei den Verbrauchern der Aktivität aus. Falls eine Ressource ausfällt, von der die Aktivität abhängig ist, befindet sich die Aktivität im Zustand CANCELED und der Zustand aller Verbraucher ändert sich in CASCADE_FAILED.

Datenknoten und Vorbedingungen

Wenn als Ausgabe einer ausgefallenen Aktivität ein Datenknoten konfiguriert ist, ändert sich der Zustand des Datenknotens in CASCADE_FAILED. Der Ausfall eines Datenknotens wird auf alle zugehörigen Vorbedingungen übertragen, deren Zustand sich in CANCELED ändert.

Ressourcen

Wenn sich die Objekte, die von einer Ressource abhängig sind, im Zustand FAILED befinden, während sich die Ressource im Zustand WAITING_ON_DEPENDENCIES befindet, dann ändert sich der Zustand der Ressource in FINISHED.

Objekte, bei denen kaskadierte Fehler aufgetreten sind, erneut ausführen

Standardmäßig wird bei der erneuten Ausführung einer Aktivität oder eines Datenknotens nur die zugehörige Ressource erneut ausgeführt. Wird jedoch das Feld failureAndRerunMode in einem Pipeline-Objekt auf cascade eingestellt, ermöglicht dies unter den folgenden Bedingungen, dass die erneute Ausführung auf einem Zielobjekt auf alle Verbraucher übertragen wird:

- Die Verbraucher des Zielobjekts befinden sich im Zustand CASCADE_FAILED.

- Für die Abhängigkeiten des Zielobjekts stehen keine Befehle für eine erneute Ausführung an.
- Die Abhängigkeiten des Zielobjekts befinden sich nicht im Zustand FAILED, CASCADE_FAILED oder CANCELED.

Bei dem Versuch, ein CASCADE_FAILED-Objekt erneut auszuführen, bei dem sich irgendwelche seiner Abhängigkeiten im Zustand FAILED, CASCADE_FAILED oder CANCELED befinden, schlägt die erneute Ausführung fehl und das Objekt kehrt wieder in den Zustand CASCADE_FAILED zurück. Zur erfolgreichen erneuten Ausführung des ausgefallenen Objekts müssen Sie den Ausfall in der Kette der Abhängigkeiten nach oben verfolgen, um die ursprüngliche Quelle des Ausfalls zu finden, und statt dessen das betreffende Objekt erneut ausführen. Wenn Sie den Befehl für eine erneute Ausführung auf einer Ressource ausgeben, versuchen Sie, alle von ihr abhängigen Ressourcen ebenfalls erneut auszuführen.

Cascade-Fehler und Backfills

Wenn Sie Cascade Failure aktivieren und über eine Pipeline verfügen, die viele Backfills erzeugt, können Pipeline-Laufzeitfehler dazu führen, dass Ressourcen schnell hintereinander erstellt und gelöscht werden, ohne dass nützliche Arbeit geleistet wird. AWS Data Pipeline versucht, Sie beim Speichern einer Pipeline mit der folgenden Warnmeldung auf diese Situation aufmerksam zu machen: *Pipeline_object_name* has 'failureAndRerunMode' field set to 'cascade' and you are about to create a backfill with scheduleStartTime *start_time*. This can result in rapid creation of pipeline objects in case of failures. Dies liegt daran, dass bei einem Kaskadenausfall nachgelagerte Aktivitäten schnell als EMR-Cluster CASCADE_FAILED und nicht mehr benötigte EC2 Ressourcen eingerichtet und heruntergefahren werden können. Wir empfehlen, dass Sie Pipelines mit kurzen Zeitbereichen testen, um die Auswirkungen dieser Situation zu begrenzen.

Syntax der Pipeline-Definitionsdatei

Die Anweisungen in diesem Abschnitt beziehen sich auf die manuelle Arbeit mit Pipeline-Definitionsdateien über die AWS Data Pipeline Befehlszeilenschnittstelle (CLI). Dies ist eine Alternative zum interaktiven Entwerfen einer Pipeline mithilfe der AWS Data Pipeline Konsole.

Sie können Pipeline-Definitionsdateien manuell mit einem beliebigen Texteditor erstellen, der das Speichern von Dateien im UTF-8-Dateiformat unterstützt, und die Dateien über die AWS Data Pipeline Befehlszeilenschnittstelle senden.

AWS Data Pipeline unterstützt auch eine Vielzahl komplexer Ausdrücke und Funktionen innerhalb von Pipelinedefinitionen. Weitere Informationen finden Sie unter [Pipeline-Ausdrücke und -Funktionen](#).

Dateistruktur

Der erste Schritt bei der Erstellung einer Pipeline besteht darin, Pipeline-Definitionsobjekte in einer Pipeline-Definitionsdatei zu verfassen. Im folgenden Beispiel wird die allgemeine Struktur einer Pipeline-Definitionsdatei veranschaulicht. Diese Datei definiert zwei Objekte, die durch '{' und '}' begrenzt und durch ein Komma getrennt werden.

Im folgenden Beispiel definiert das erste Objekt zwei Namen-Wert-Paare, die als Felder bezeichnet werden. Das zweite Objekt definiert drei Felder.

```
{  
  "objects" : [  
    {  
      "name1" : "value1",  
      "name2" : "value2"  
    },  
    {  
      "name1" : "value3",  
      "name3" : "value4",  
      "name4" : "value5"  
    }  
  ]  
}
```

Beim Erstellen einer Pipeline-Definitionsdatei müssen Sie die erforderlichen Typen von Pipeline-Objekten auswählen, sie der Pipeline-Definitionsdatei hinzufügen und dann die entsprechenden Felder hinzufügen. Weitere Informationen zu Pipeline-Objekten finden Sie unter [Pipeline-Objektreferenz](#).

So können Sie beispielsweise ein Pipeline-Definitionsobjekt für einen Eingabedatenknoten und ein anderes für den Ausgabedatenknoten erstellen. Erstellen Sie dann ein weiteres Pipeline-Definitionsobjekt für eine Aktivität, z. B. die Verarbeitung der Eingabedaten mit Amazon EMR.

Pipeline-Felder

Wenn Sie wissen, welche Objekttypen in Ihre Pipeline-Definitionsdatei aufzunehmen sind, fügen Sie die betreffenden Felder zur Definition der einzelnen Pipeline-Objekte hinzu. Feldnamen stehen in

Anführungszeichen und sind wie im folgenden Beispiel veranschaulicht durch ein Leerzeichen, einen Doppelpunkt und ein Leerzeichen von den Feldwerten getrennt.

```
"name" : "value"
```

Bei dem Feldwert kann es sich um eine Zeichenfolge, einen Verweis auf ein anderes Objekt, einen Funktionsaufruf, einen Ausdruck oder eine geordnete Liste beliebiger der oben genannten Typen handeln. Weitere Informationen über die Arten von Daten, die für Feldwerte verwendet werden können, finden Sie unter [Einfache Datentypen](#). Weitere Informationen zu Funktionen, die zur Auswertung von Feldwerten verwendet werden können, finden Sie unter [Ausdrucksauswertung](#).

Felder sind auf 2048 Zeichen begrenzt. Da die Größe von Objekten 20 KB betragen kann, können einem Objekt nicht viele große Felder hinzugefügt werden.

Jedes Pipeline-Objekt muss die folgenden Felder enthalten: `id` und `type`, wie im folgenden Beispiel dargestellt. Je nach Objekttyp werden möglicherweise noch andere Felder benötigt. Wählen Sie einen Wert für `id` aus, der für Sie bedeutsam und in der Pipeline-Definition eindeutig ist. Der Wert für `type` gibt den Typ des Objekts an. Geben Sie einen der unterstützten Objekttypen für die Pipeline-Definition ein, die im Thema [Pipeline-Objektreferenz](#) aufgelistet werden.

```
{
  "id": "MyCopyToS3",
  "type": "CopyActivity"
}
```

Weitere Informationen zu den erforderlichen und optionalen Feldern eines jeden Objekts finden Sie in der Dokumentation für das Objekt.

Wenn Sie Felder aus einem Objekt in ein anderes Objekt einschließen möchten, verwenden Sie das Feld `parent` mit einem Verweis auf das Objekt. Beispiel: Objekt „B“ enthält seine Felder, „B1“ und „B2“, sowie die Felder von Objekt „A“, „A1“ und „A2“.

```
{
  "id" : "A",
  "A1" : "value",
  "A2" : "value"
},
{
  "id" : "B",
```

```
"parent" : {"ref" : "A"},  
"B1" : "value",  
"B2" : "value"  
}
```

Sie können gemeinsame Felder in einem Objekt mit der ID „Default“ definieren. Diese Felder werden automatisch in die Pipeline-Definitionsdatei eines jeden Objekts eingeschlossen, sofern für das Feld `parent` nicht ausdrücklich ein Verweis auf ein anderes Objekt festgelegt ist.

```
{  
  "id" : "Default",  
  "onFail" : {"ref" : "FailureNotification"},  
  "maximumRetries" : "3",  
  "workerGroup" : "myWorkerGroup"  
}
```

Benutzerdefinierte Felder

Sie können für Ihre Pipeline-Komponenten benutzerdefinierte Felder erstellen und mit Ausdrücken auf sie verweisen. Das folgende Beispiel zeigt ein benutzerdefiniertes Feld, das benannt `myCustomField` und einem DataNode S3-Objekt `my_customFieldReference` hinzugefügt wurde:

```
{  
  "id": "S3DataInput",  
  "type": "S3DataNode",  
  "schedule": {"ref": "TheSchedule"},  
  "filePath": "s3://bucket_name",  
  "myCustomField": "This is a custom value in a custom field.",  
  "my_customFieldReference": {"ref": "AnotherPipelineComponent"}  
},
```

Der Name eines benutzerdefinierten Feldes muss als Präfix das Wort „`my`“ in Kleinbuchstaben gefolgt von einem Großbuchstaben und einem Unterstrich besitzen. Darüber hinaus kann ein benutzerdefiniertes Feld wie im voranstehenden `myCustomField`-Beispiel ein Zeichenfolgenwert oder wie im voranstehenden `my_customFieldReference`-Beispiel ein Verweis auf eine andere Pipeline-Komponente sein.

Note

Prüft bei benutzerdefinierten Feldern AWS Data Pipeline nur auf gültige Verweise auf andere Pipeline-Komponenten, nicht auf Zeichenkettenwerte für benutzerdefinierte Felder, die Sie hinzufügen.

Arbeiten mit der API

Note

Wenn Sie keine Programme schreiben, die mit interagieren AWS Data Pipeline, müssen Sie keine der AWS installieren SDKs. Sie können über die Konsole oder die Befehlszeile Pipelines erstellen und ausführen. Weitere Informationen finden Sie unter [Einrichtung für AWS Data Pipeline](#)

Der einfachste Weg, Anwendungen zu schreiben, die mit einem benutzerdefinierten Task Runner interagieren AWS Data Pipeline oder ihn implementieren, ist die Verwendung eines der AWS SDKs. AWS SDKs bietet Funktionen, die das Aufrufen des Webservices APIs von Ihrer bevorzugten Programmierumgebung aus vereinfachen. Weitere Informationen finden Sie unter [Installieren des AWS-SDKs](#).

Installieren des AWS-SDKs

Sie AWS SDKs bieten Funktionen, die die API umschließen und sich um viele Verbindungsdetails kümmern, z. B. um die Berechnung von Signaturen, die Bearbeitung von Wiederholungsversuchen von Anfragen und die Fehlerbehandlung. Sie enthalten SDKs auch Beispielcode, Tutorials und andere Ressourcen, die Ihnen den Einstieg in das Schreiben von Anwendungen erleichtern, die aufrufen AWS. Das Aufrufen der Wrapper-Funktionen in einem SDK kann das Schreiben einer AWS Anwendung erheblich vereinfachen. Weitere Informationen zum Herunterladen und Verwenden von finden Sie AWS SDKs unter [Beispielcode und Bibliotheken](#).

AWS Data Pipeline Support ist SDKs für die folgenden Plattformen verfügbar:

- [AWS SDK für Java](#)
- [AWS SDK for Node.js](#)

- [AWS SDK für PHP](#)
- [AWS SDK für Python \(Boto\)](#)
- [AWS SDK für Ruby](#)
- [AWS SDK für .NET](#)

Eine HTTP-Anfrage stellen an AWS Data Pipeline

Eine vollständige Beschreibung der programmatischen Objekte finden Sie in AWS Data Pipeline der [AWS Data Pipeline API-Referenz](#).

Wenn Sie keine der AWS SDKs verwenden möchten, können Sie mithilfe der POST-Anforderungsmethode AWS Data Pipeline Operationen über HTTP ausführen. Bei der POST-Methode müssen Sie den Vorgang im Header der Anforderung festlegen und im Anforderungstext die Daten für den Vorgang im JSON-Format angeben.

Inhalt des HTTP-Headers

AWS Data Pipeline erfordert die folgenden Informationen im Header einer HTTP-Anfrage:

- host Der AWS Data Pipeline Endpunkt.
Weitere Informationen zu Endpunkten finden Sie unter [Regionen und Endpunkte](#).
- x-amz-date Sie müssen den Zeitstempel entweder im HTTP-Date-Header oder im x-amz-date AWS-Header angeben. (Einige HTTP-Client-Bibliotheken lassen den Datums-Header nicht zu.) Wenn ein x-amz-date Header vorhanden ist, ignoriert das System alle Date-Header während der Anforderungsaufentifizierung.

Das Datum muss in einem der folgenden drei Formate angegeben werden, wie in HTTP/1.1 RFC festgelegt:

- Sun, 06 Nov 1994 08:49:37 GMT (RFC 822, aktualisiert durch RFC 1123)
- Sunday, 06-Nov-94 08:49:37 GMT (RFC 850, abgelöst durch RFC 1036)
- Sun Nov 6 08:49:37 1994 (ANSI C asctime()-Format)
- Authorization Der Satz an Autorisierungsparametern, mit denen AWS die Gültigkeit und Authentizität der Anforderung sicherstellt. Weitere Informationen zum Aufbau dieses Headers finden Sie unter [Signature Version 4-Signaturprozess](#).
- x-amz-target Der Zieldienst der Anforderung und des Vorgangs für die Daten, in folgendem Format: <>service><>_<>API version><>. <>operationName><>

Beispiel: DataPipeline_20121129.ActivatePipeline

- content-type Gibt JSON und die Version an. Beispiel: Content-Type: application/x-amz-json-1.0

Nachfolgend finden Sie einen Beispiel-Header für eine HTTP-Anforderung zum Aktivieren einer Pipeline.

```
POST / HTTP/1.1
host: https://datapipeline.us-east-1.amazonaws.com
x-amz-date: Mon, 12 Nov 2012 17:49:52 GMT
x-amz-target: DataPipeline_20121129.ActivatePipeline
Authorization: AuthParams
Content-Type: application/x-amz-json-1.1
Content-Length: 39
Connection: Keep-Alive
```

HTTP-Textinhalt

Der Textkörper einer HTTP-Anforderung enthält die Daten für den Vorgang, der im Header der HTTP-Anforderung festgelegt ist. Die Daten müssen gemäß dem JSON-Datenschema für jede AWS Data Pipeline API formatiert werden. Das AWS Data Pipeline JSON-Datenschema definiert die Arten von Daten und Parametern (wie Vergleichsoperatoren und Aufzählungskonstanten), die für jede Operation verfügbar sind.

Format des Textkörpers einer HTTP-Anforderung

Verwenden Sie das JSON-Datenformat zur gleichzeitigen Übermittlung von Datenwerten und -strukturen. Elemente können mit der Klammerschreibweise innerhalb anderer Elemente verschachtelt werden. Das folgende Beispiel zeigt eine Anforderung für das Erstellen einer Pipeline-Definition, die aus drei Objekten und ihren entsprechenden Slots besteht.

```
{
  "pipelineId": "df-00627471SOVYZEXAMPLE",
  "pipelineObjects":
  [
    {"id": "Default",
     "name": "Default",
```

```
"slots":  
  [  
    {"key": "workerGroup",  
     "stringValue": "MyWorkerGroup"}  
  ]  
,  
 {"id": "Schedule",  
  "name": "Schedule",  
  "slots":  
    [  
      {"key": "startDateTime",  
       "stringValue": "2012-09-25T17:00:00"},  
      {"key": "type",  
       "stringValue": "Schedule"},  
      {"key": "period",  
       "stringValue": "1 hour"},  
      {"key": "endDateTime",  
       "stringValue": "2012-09-25T18:00:00"}  
    ]  
,  
 {"id": "SayHello",  
  "name": "SayHello",  
  "slots":  
    [  
      {"key": "type",  
       "stringValue": "ShellCommandActivity"},  
      {"key": "command",  
       "stringValue": "echo hello"},  
      {"key": "parent",  
       "refValue": "Default"},  
      {"key": "schedule",  
       "refValue": "Schedule"}  
    ]  
,  
 }]
```

Handhaben der HTTP-Antwort

Nachfolgend finden Sie einige wichtige Header in der HTTP-Antwort und Informationen dazu, wie Sie diese Header in Ihrer Anwendung behandeln sollten:

- HTTP/1.1 — Auf diesen Header folgt ein Statuscode. Ein Code-Wert von 200 gibt an, dass ein Vorgang erfolgreich war. Jeder andere Wert weist auf einen Fehler hin.
- x-amzn- RequestId — Dieser Header enthält eine Anforderungs-ID, die Sie verwenden können, wenn Sie Probleme mit einer Anfrage beheben müssen. AWS Data Pipeline Ein Beispiel für eine Anforderungs-ID ist QH8 DNOU9 K2 07N97 Q9ASUAAJG. FNA2 GDLL8 OBVV4 KQNSO5 AEMVJF66
- x-amz-crc32 —AWS Data Pipeline berechnet eine CRC32 Prüfsumme der HTTP-Nutzlast und gibt diese Prüfsumme im 32-Header zurück. x-amz-crc Wir empfehlen, dass Sie Ihre eigene CRC32 Prüfsumme auf der Client-Seite berechnen und sie mit dem x-amz-crc 32-Header vergleichen. Wenn die Prüfsummen nicht übereinstimmen, kann dies darauf hinweisen, dass die Daten während der Übertragung beschädigt wurden. Wenn dies der Fall ist, sollten Sie Ihre Anforderung erneut übermitteln.

AWS-SDK-Benutzer müssen diese Überprüfung nicht manuell durchführen, da sie die Prüfsumme jeder Antwort von Amazon DynamoDB SDKs berechnen und es automatisch erneut versuchen, wenn eine Nichtübereinstimmung festgestellt wird.

Beispiel für eine JSON-Anfrage und -Antwort AWS Data Pipeline

Die folgenden Beispiele zeigen eine Anforderung zum Erstellen einer neuen Pipeline. Anschließend wird die AWS Data Pipeline Antwort einschließlich der Pipeline-ID der neu erstellten Pipeline angezeigt.

HTTP-POST-Anforderung

```
POST / HTTP/1.1
host: https://datapipeline.us-east-1.amazonaws.com
x-amz-date: Mon, 12 Nov 2012 17:49:52 GMT
x-amz-target: DataPipeline_20121129.CreatePipeline
Authorization: AuthParams
Content-Type: application/x-amz-json-1.1
Content-Length: 50
Connection: Keep-Alive

{"name": "MyPipeline",
 "uniqueId": "12345ABCDEFG"}
```

AWS Data Pipeline Antwort

```
HTTP/1.1 200
x-amzn-RequestId: b16911ce-0774-11e2-af6f-6bc7a6be60d9
x-amz-crc32: 2215946753
Content-Type: application/x-amz-json-1.0
Content-Length: 2
Date: Mon, 16 Jan 2012 17:50:53 GMT

{"pipelineId": "df-00627471S0VYZEXAMPLE"}
```

Sicherheit in AWS Data Pipeline

Cloud-Sicherheit AWS hat höchste Priorität. Als AWS Kunde profitieren Sie von Rechenzentren und Netzwerkarchitekturen, die darauf ausgelegt sind, die Anforderungen der sicherheitssensibelsten Unternehmen zu erfüllen.

Sicherheit ist eine gemeinsame AWS Verantwortung von Ihnen und Ihnen. Das [Modell der geteilten Verantwortung](#) beschreibt dies als Sicherheit der Cloud und Sicherheit in der Cloud:

- Sicherheit der Cloud — AWS ist verantwortlich für den Schutz der Infrastruktur, die AWS Dienste in der AWS Cloud ausführt. AWS bietet Ihnen auch Dienste, die Sie sicher nutzen können. Externe Prüfer testen und verifizieren regelmäßig die Wirksamkeit unserer Sicherheitsmaßnahmen im Rahmen der [AWS](#). Weitere Informationen zu den geltenden Compliance-Programmen finden Sie unter [AWS-Services in Umfang nach Compliance-Programm](#). AWS Data Pipeline
- Sicherheit in der Cloud — Ihre Verantwortung richtet sich nach dem AWS Service, den Sie nutzen. Sie sind auch für andere Faktoren verantwortlich, etwa für die Vertraulichkeit Ihrer Daten, für die Anforderungen Ihres Unternehmens und für die geltenden Gesetze und Vorschriften.

Diese Dokumentation hilft Ihnen zu verstehen, wie Sie das Modell der gemeinsamen Verantwortung bei der Nutzung anwenden können AWS Data Pipeline. In den folgenden Themen erfahren Sie, wie Sie die Konfiguration vornehmen AWS Data Pipeline , um Ihre Sicherheits- und Compliance-Ziele zu erreichen. Sie lernen auch, wie Sie andere AWS-Services nutzen können, die Sie bei der Überwachung und Sicherung Ihrer AWS Data Pipeline Ressourcen unterstützen.

Topics

- [Datenschutz in AWS Data Pipeline](#)
- [Identity and Access Management für AWS Data Pipeline](#)
- [Einloggen und Überwachen AWS Data Pipeline](#)
- [Reaktion auf Vorfälle in AWS Data Pipeline](#)
- [Überprüfung der Einhaltung von Vorschriften für AWS Data Pipeline](#)
- [Resilienz in AWS Data Pipeline](#)
- [Sicherheit der Infrastruktur in AWS Data Pipeline](#)
- [Konfiguration und Schwachstellenanalyse in AWS Data Pipeline](#)

Datenschutz in AWS Data Pipeline

Das [Modell der AWS gemeinsamen Verantwortung](#) und geteilter Verantwortung gilt für den Datenschutz in AWS Data Pipeline. Wie in diesem Modell beschrieben, AWS ist verantwortlich für den Schutz der globalen Infrastruktur, auf der alle Systeme laufen AWS Cloud. Sie sind dafür verantwortlich, die Kontrolle über Ihre in dieser Infrastruktur gehosteten Inhalte zu behalten. Dieser Inhalt enthält die Sicherheitskonfigurations- und Verwaltungsaufgaben für die von Ihnen verwendeten AWS-Services . Weitere Informationen zum Datenschutz finden Sie unter [Häufig gestellte Fragen zum Datenschutz](#). Informationen zum Datenschutz in Europa finden Sie im Blog-Beitrag [AWS -Modell der geteilten Verantwortung und in der DSGVO](#) im AWS -Sicherheitsblog.

Aus Datenschutzgründen empfehlen wir, dass Sie AWS-Konto Anmeldeinformationen schützen und einzelne Benutzer mit AWS IAM Identity Center oder AWS Identity and Access Management (IAM) einrichten. So erhält jeder Benutzer nur die Berechtigungen, die zum Durchführen seiner Aufgaben erforderlich sind. Außerdem empfehlen wir, die Daten mit folgenden Methoden schützen:

- Verwenden Sie für jedes Konto die Multi-Faktor-Authentifizierung (MFA).
- Verwenden Sie SSL/TLS, um mit Ressourcen zu kommunizieren. AWS Wir empfehlen TLS 1.2 oder höher.
- Richten Sie die API und die Protokollierung von Benutzeraktivitäten mit ein. AWS CloudTrail
- Verwenden Sie AWS Verschlüsselungslösungen zusammen mit allen darin enthaltenen Standardsicherheitskontrollen AWS-Services.
- Verwenden Sie erweiterte verwaltete Sicherheitsservices wie Amazon Macie, die dabei helfen, in Amazon S3 gespeicherte persönliche Daten zu erkennen und zu schützen.
- Wenn Sie für den Zugriff AWS über eine Befehlszeilenschnittstelle oder eine API FIPS 140-2- validierte kryptografische Module benötigen, verwenden Sie einen FIPS-Endpunkt. Weitere Informationen über verfügbare FIPS-Endpunkte finden Sie unter [Federal Information Processing Standard \(FIPS\) 140-2](#).
- AWS Data Pipeline unterstützt IMDSv2 Amazon EMR und EC2 Amazon-Ressourcen. Verwenden Sie zur Verwendung IMDSv2 mit Amazon EMR die Versionen 5.23.1, 5.27.1 oder 5.32 oder höher oder Version 6.2 oder höher. Weitere Informationen finden [Sie unter Metadaten-Serviceanfragen an EC2 Amazon-Instances konfigurieren](#) und [verwenden IMDSv2](#).

Wir empfehlen dringend, in Freitextfeldern, z. B. im Feld Name, keine vertraulichen oder sensiblen Informationen wie die E-Mail-Adressen Ihrer Kunden einzugeben. Dies gilt auch, wenn Sie mit der Konsole, der AWS-Services API AWS Data Pipeline oder auf andere Weise arbeiten oder

diese verwenden AWS SDKs. AWS CLI Alle Daten, die Sie in Tags oder Freitextfelder eingeben, die für Namen verwendet werden, können für Abrechnungs- oder Diagnoseprotokolle verwendet werden. Wenn Sie eine URL für einen externen Server bereitstellen, empfehlen wir dringend, keine Anmeldeinformationen zur Validierung Ihrer Anforderung an den betreffenden Server in die URL einzuschließen.

Identity and Access Management für AWS Data Pipeline

Mit Ihren Sicherheitsanmeldeinformationen identifizieren Sie sich bei den Services in AWS und erhalten Berechtigungen zur Verwendung von AWS-Ressourcen, wie etwa Ihre Pipelines. Sie können die Funktionen von AWS Data Pipeline und AWS Identity and Access Management (IAM) verwenden, um anderen Benutzern den Zugriff auf Ihre AWS Data Pipeline Ressourcen zu ermöglichen AWS Data Pipeline , ohne Ihre Sicherheitsanmeldedaten weitergeben zu müssen.

In Unternehmen können bestimmte Pipelines für mehrere Mitarbeiter freigegeben werden, damit diese gemeinsam damit arbeiten können. Allerdings sollten in diesem Fall Maßnahmen wie die folgenden ergriffen werden:

- Steuern Sie, welche Benutzer auf bestimmte Pipelines zugreifen können
- Schützen einer Produktions-Pipeline vor versehentlichen Änderungen
- Erlauben des Lesezugriffs von Auditoren auf Pipelines und gleichzeitiges Verhindern von Änderungen

AWS Data Pipeline ist in AWS Identity and Access Management (IAM) integriert, das eine Vielzahl von Funktionen bietet:

- Erstellen Sie Benutzer und Gruppen in Ihrem AWS-Konto.
- Teilen Sie Ihre AWS Ressourcen ganz einfach zwischen den Benutzern in Ihrem AWS-Konto.
- Weisen Sie jedem Benutzer eindeutige Sicherheitsanmeldedaten zu.
- Kontrollieren Sie den Zugriff jedes Benutzers auf Dienste und Ressourcen.
- Holen Sie sich eine einzige Rechnung für alle Benutzer in Ihrem AWS-Konto.

Durch die Verwendung von IAM mit können Sie steuern AWS Data Pipeline, ob Benutzer in Ihrer Organisation eine Aufgabe mithilfe bestimmter API-Aktionen ausführen können und ob sie bestimmte AWS-Ressourcen verwenden können. Sie können IAM-Richtlinien verwenden, die auf Pipeline-

Tags und Worker-Gruppen basieren, um Ihre Pipelines mit anderen Benutzern zu teilen und deren Zugriffsebene zu kontrollieren.

Inhalt

- [IAM-Richtlinien für AWS Data Pipeline](#)
- [Beispielrichtlinien für AWS Data Pipeline](#)
- [IAM-Rollen für AWS Data Pipeline](#)

IAM-Richtlinien für AWS Data Pipeline

Standardmäßig sind IAM-Entitäten nicht berechtigt, AWS-Ressourcen zu erstellen oder zu ändern. Damit IAM-Entitäten Ressourcen erstellen oder ändern und Aufgaben ausführen können, müssen Sie IAM-Richtlinien erstellen, die IAM-Entitäten die Erlaubnis erteilen, die spezifischen Ressourcen und API-Aktionen zu verwenden, die sie benötigen, und diese Richtlinien dann den IAM-Entitäten zuordnen, für die diese Berechtigungen erforderlich sind.

Wenn Sie einem Benutzer oder einer Benutzergruppe eine Richtlinie zuordnen, wird den Benutzern die Ausführung der angegebenen Aufgaben für die angegebenen Ressourcen gestattet oder verweigert. Allgemeine Informationen zu IAM-Richtlinien finden Sie unter [Berechtigungen und Richtlinien](#) im IAM-Benutzerhandbuch. Weitere Informationen zum Verwalten und Erstellen von benutzerdefinierten IAM-Richtlinien finden Sie unter [Verwalten von IAM-Richtlinien](#).

Inhalt

- [Richtliniensyntax](#)
- [Steuern des Zugriffs auf Pipelines mithilfe von Tags](#)
- [Steuern des Zugriffs auf Pipelines mithilfe von Worker-Gruppen](#)

Richtliniensyntax

Eine IAM-Richtlinie ist ein JSON-Dokument, das eine oder mehrere Anweisungen enthält. Jede Anweisung ist folgendermaßen strukturiert:

```
{  
  "Statement": [ {  
    "Effect": "effect",  
    "Action": "action",  
    "Resource": "*" } ]}
```

```
"Condition":{  
    "condition":{  
        "key":"value"  
    }  
}  
}  
]  
}
```

Eine Anweisung in einer Richtlinie besteht aus folgenden Elementen:

- Effect: Der effect-Wert kann Allow oder Deny lauten. Standardmäßig sind IAM-Entitäten nicht berechtigt, Ressourcen und API-Aktionen zu verwenden, sodass alle Anfragen abgelehnt werden. Dieser Standardwert kann durch eine explizite Zugriffserlaubnis überschrieben werden. Eine explizite Zugriffsverweigerung überschreibt jedwede Zugriffserlaubnis.
- Action: Mit action wird die API-Aktion spezifiziert, für die Sie Berechtigungen erteilen oder verweigern. Eine Liste der Aktionen für finden Sie AWS Data Pipeline unter [Aktionen](#) in der AWS Data Pipeline API-Referenz.
- Resource: Die von einer Aktion betroffene Ressource. Der einzige hier zulässige Wert lautet "*".
- Condition: Bedingungen sind optional. Mit ihrer Hilfe können Sie bestimmen, wann Ihre Richtlinie wirksam wird.

AWS Data Pipeline implementiert die AWS-weiten Kontextschlüssel (siehe [Verfügbare Schlüssel für Bedingungen](#)) sowie die folgenden dienstspezifischen Schlüssel.

- `datipeline:PipelineCreator`— Um dem Benutzer, der die Pipeline erstellt hat, Zugriff zu gewähren. Ein Beispiel hierzu finden Sie unter [Gewähren des vollen Zugriffs für Pipeline-Eigentümer](#).
- `datipeline:Tag`— Um Zugriff auf der Grundlage von Pipeline-Tagging zu gewähren. Weitere Informationen finden Sie unter [Steuern des Zugriffs auf Pipelines mithilfe von Tags](#).
- `datipeline:workerGroup`— Um Zugriff auf der Grundlage des Namens der Worker-Gruppe zu gewähren. Weitere Informationen finden Sie unter [Steuern des Zugriffs auf Pipelines mithilfe von Worker-Gruppen](#).

Steuern des Zugriffs auf Pipelines mithilfe von Tags

Sie können IAM-Richtlinien erstellen, die auf die Tags für Ihre Pipeline verweisen. Dadurch lässt sich mithilfe von Pipeline-Tags Folgendes durchführen:

- Gewähren des Lesezugriffs auf eine Pipeline
- read/write Zugriff auf eine Pipeline gewähren
- Blockieren des Zugriffs auf eine Pipeline

Nehmen wir an, dass es in einem Unternehmen zwei Pipeline-Umgebungen (Produktion und Entwicklung) und für jede Umgebung eine IAM-Gruppe gibt. Für Pipelines in der Produktionsumgebung gewährt der Manager Benutzern in der IAM-Produktionsgruppe read/write Zugriff, Benutzern in der Entwickler-IAM-Gruppe jedoch nur Lesezugriff. Für Pipelines in der Entwicklungsumgebung gewährt der Manager read/write Zugriff sowohl auf die Produktions- als auch auf die Entwickler-IAM-Gruppe.

Um dieses Szenario zu erreichen, kennzeichnet der Manager die Produktionspipelines mit dem Tag „environment=production“ und fügt der Entwickler-IAM-Gruppe die folgende Richtlinie hinzu. Die erste Anweisung gewährt Lesezugriff auf alle Pipelines. Die zweite Anweisung gewährt read/write Zugriff auf Pipelines, die nicht über das Tag „environment=production“ verfügen.

JSON

```
{  
    "Version": "2012-10-17",  
    "Statement": [  
        {  
            "Effect": "Allow",  
            "Action": [  
                "datapipeline:Describe*",  
                "datapipeline>ListPipelines",  
                "datapipeline:GetPipelineDefinition",  
                "datapipeline:QueryObjects"  
            ],  
            "Resource": "*"  
        },  
        {  
            "Effect": "Allow",  
            "Action": "datapipeline:*",  
            "Resource": "*",  
            "Condition": {  
                "StringNotEquals": {"datapipeline:Tag/environment": "production"}  
            }  
        }  
    ]}
```

{

Darüber hinaus weist der Manager der IAM-Produktionsgruppe die folgende Richtlinie zu. Diese Anweisung gewährt vollständigen Zugriff auf alle Pipelines.

JSON

```
{  
    "Version": "2012-10-17",  
    "Statement": [  
        {  
            "Effect": "Allow",  
            "Action": "datapipeline:*",  
            "Resource": "*"  
        }  
    ]  
}
```

Weitere Beispiele finden Sie unter [Gewähren des Lesezugriffs für Benutzer basierend auf einem Tag](#) und [Gewähren des vollständigen Zugriffs für Benutzer basierend auf einem Tag](#).

Steuern des Zugriffs auf Pipelines mithilfe von Worker-Gruppen

Sie können IAM-Richtlinien erstellen, die auf Namen von Workergruppen verweisen.

Nehmen wir an, dass es in einem Unternehmen zwei Pipeline-Umgebungen (Produktion und Entwicklung) und für jede Umgebung eine IAM-Gruppe gibt. Außerdem sind drei Datenbankserver mit Task Runner-Anwendungen für die Produktions-, die Vorproduktions- und die Entwicklungsumgebung vorhanden. Der Manager möchte sicherstellen, dass Benutzer in der IAM-Produktionsgruppe Pipelines erstellen können, die Aufgaben an Produktionsressourcen weiterleiten, und dass Benutzer in der IAM-Entwicklungsgruppe Pipelines erstellen können, die Aufgaben sowohl an Vorproduktions- als auch an Entwicklerressourcen weiterleiten.

Um dies zu erreichen, installiert er Task Runner auf den Produktionsressourcen mit den Anmeldeinformationen der Produktionsgruppe und weist `workerGroup` den Wert "prodresource" zu. Außerdem installiert der Manager Task Runner auf den Entwicklungsressourcen mit den Anmeldeinformationen der Entwicklungsgruppe und weist `workerGroup` die Werte "pre-production" und "development" zu. Der Manager weist der Entwickler-IAM-Gruppe die folgende Richtlinie zu, um

den Zugriff auf „prodresource“ -Ressourcen zu blockieren. Die erste Anweisung gewährt Lesezugriff auf alle Pipelines. Die zweite Anweisung gewährt read/write Zugriff auf Pipelines, wenn der Name der Worker-Gruppe das Präfix „dev“ oder „pre-prod“ hat.

Darüber hinaus weist der Manager der IAM-Produktionsgruppe die folgende Richtlinie zu, um Zugriff auf „prodresource“ -Ressourcen zu gewähren. Die erste Anweisung gewährt Lesezugriff auf alle Pipelines. Die zweite Anweisung gewährt read/write Zugriff, wenn der Name der Worker-Gruppe das Präfix „prod“ hat.

Beispielrichtlinien für AWS Data Pipeline

Die folgenden Beispiele zeigen, wie Sie Benutzern vollständigen oder eingeschränkten Zugriff auf Pipelines gewähren.

Inhalt

- [Beispiel 1: Gewähren des Lesezugriffs für Benutzer basierend auf einem Tag](#)
- [Beispiel 2: Gewähren des vollständigen Zugriffs für Benutzer basierend auf einem Tag](#)
- [Beispiel 3: Gewähren des vollen Zugriffs für Pipeline-Eigentümer](#)
- [Beispiel 4: Gewähren Sie Benutzern Zugriff auf die Konsole AWS Data Pipeline](#)

Beispiel 1: Gewähren des Lesezugriffs für Benutzer basierend auf einem Tag

Die folgende Richtlinie ermöglicht es Benutzern, die schreibgeschützten AWS Data Pipeline API-Aktionen zu verwenden, jedoch nur mit Pipelines, die das Tag „environment=production“ haben.

Die ListPipelines API-Aktion unterstützt keine Tag-basierte Autorisierung.

JSON

```
{  
    "Version": "2012-10-17",  
    "Statement": [  
        {  
            "Effect": "Allow",  
            "Action": [  
                "datapipeline:Describe*",  
                "datapipeline:GetPipelineDefinition",  
                "datapipeline:ValidatePipelineDefinition",  
                "datapipeline:PutPipelineDefinition",  
                "datapipeline:DeletePipeline"  
            ]  
        }  
    ]  
}
```

```
        "datapipeline:QueryObjects"
    ],
    "Resource": [
        "*"
    ],
    "Condition": {
        "StringEquals": {
            "datapipeline:Tag/environment": "production"
        }
    }
}
]
```

Beispiel 2: Gewähren des vollständigen Zugriffs für Benutzer basierend auf einem Tag

Die folgende Richtlinie ermöglicht es Benutzern, alle AWS Data Pipeline API-Aktionen zu verwenden, mit Ausnahme von ListPipelines, aber nur mit Pipelines, die das Tag „environment=test“ haben.

JSON

```
{
    "Version": "2012-10-17",
    "Statement": [
        {
            "Effect": "Allow",
            "Action": [
                "datapipeline:*"
            ],
            "Resource": [
                "*"
            ],
            "Condition": {
                "StringEquals": {
                    "datapipeline:Tag/environment": "test"
                }
            }
        }
    ]
}
```

Beispiel 3: Gewähren des vollen Zugriffs für Pipeline-Eigentümer

Die folgende Richtlinie ermöglicht es Benutzern, alle AWS Data Pipeline API-Aktionen zu verwenden, jedoch nur mit ihren eigenen Pipelines.

Beispiel 4: Gewähren Sie Benutzern Zugriff auf die Konsole AWS Data Pipeline

Die folgende Richtlinie gestattet Benutzern, mit der AWS Data Pipeline -Konsole eine Pipeline zu erstellen und zu verwalten.

Diese Richtlinie umfasst die Aktion für PassRole Berechtigungen für bestimmte Ressourcen, die an roleARN die jeweiligen AWS Data Pipeline Anforderungen gebunden sind. Weitere Informationen zur identitätsbasierten (IAM) PassRole -Berechtigung finden Sie im Blogbeitrag [Granting Permission to Launch EC2 Instances with IAM Roles \(Permission\)](#). PassRole

JSON

```
{  
    "Version": "2012-10-17",  
    "Statement": [ {  
        "Action": [  
            "cloudwatch:*",  
            "datapipeline:*",  
            "dynamodb:DescribeTable",  
            "elasticmapreduce:AddJobFlowSteps",  
            "elasticmapreduce>ListInstance*",  
            "iam:AddRoleToInstanceProfile",  
            "iam>CreateInstanceProfile",  
            "iam:GetInstanceProfile",  
            "iam:GetRole",  
            "iam:GetRolePolicy",  
            "iam>ListInstanceProfiles",  
            "iam>ListInstanceProfilesForRole",  
            "iam>ListRoles",  
            "rds:DescribeDBInstances",  
            "rds:DescribeDBSecurityGroups",  
            "redshift:DescribeClusters",  
            "redshift:DescribeClusterSecurityGroups",  
            "s3>List*",  
            "sns>ListTopics"  
        ],  
        "Effect": "Allow",  
        "Resource": "  
    }]
```

```
"Resource": [
    "*"
],
},
{
    "Action": "iam:PassRole",
    "Effect": "Allow",
    "Resource": [
        "arn:aws:iam::*:role/DataPipelineDefaultResourceRole",
        "arn:aws:iam::*:role/DataPipelineDefaultRole"
    ]
}
```

IAM-Rollen für AWS Data Pipeline

AWS Data Pipeline verwendet AWS Identity and Access Management Rollen. Die den IAM-Rollen zugewiesenen Berechtigungsrichtlinien bestimmen, welche Aktionen AWS Data Pipeline und Ihre Anwendungen ausführen können und auf welche AWS Ressourcen sie zugreifen können. Weitere Informationen finden Sie unter [IAM-Rollen](#) im IAM-Benutzerhandbuch.

AWS Data Pipeline erfordert zwei IAM-Rollen:

- Die Pipeline-Rolle steuert AWS Data Pipeline den Zugriff auf Ihre AWS-Ressourcen. In Pipeline-Objektdefinitionen gibt das `role` Feld diese Rolle an.
- Die EC2 Instance-Rolle steuert den Zugriff, den Anwendungen, die auf EC2 Instances ausgeführt werden, einschließlich der EC2 Instances in Amazon EMR-Clustern, auf AWS Ressourcen haben. In Pipeline-Objektdefinitionen gibt das `resourceRole` Feld diese Rolle an.

Important

Wenn Sie vor dem 3. Oktober 2022 mithilfe der AWS Data Pipeline Konsole mit Standardrollen eine Pipeline AWS Data Pipeline erstellt haben, haben Sie die `DataPipelineDefaultRole` für Sie erstellt und die `AWSDataPipelineRole` verwaltete Richtlinie an die Rolle angehängt. Seit dem 3. Oktober 2022 ist die `AWSDataPipelineRole` verwaltete Richtlinie veraltet und die Pipeline-Rolle muss für eine Pipeline angegeben werden, wenn Sie die Konsole verwenden.

Wir empfehlen Ihnen, die vorhandenen Pipelines zu überprüfen und festzustellen, ob die mit der Pipeline verknüpft `DataPipelineDefaultRole` ist und ob die mit dieser `AWSDataPipelineRole` Rolle verknüpft ist. Falls ja, überprüfen Sie den Zugriff, den diese Richtlinie gewährt, um sicherzustellen, dass er Ihren Sicherheitsanforderungen entspricht. Fügen Sie nach Bedarf die mit dieser Rolle verknüpften Richtlinien und Richtlinienerklärungen hinzu, aktualisieren oder ersetzen Sie sie. Alternativ können Sie eine Pipeline aktualisieren, um eine Rolle zu verwenden, die Sie mit unterschiedlichen Berechtigungsrichtlinien erstellt haben.

Beispiel für Berechtigungsrichtlinien für AWS Data Pipeline Rollen

Jeder Rolle sind eine oder mehrere Berechtigungsrichtlinien zugeordnet, die festlegen, auf welche AWS Ressourcen die Rolle zugreifen kann und welche Aktionen die Rolle ausführen kann. Dieses Thema enthält ein Beispiel für eine Berechtigungsrichtlinie für die Pipeline-Rolle. Es enthält auch den Inhalt der `AmazonEC2RoleforDataPipelineRole`, der verwalteten Richtlinie für die EC2 Standard-Instanzrolle, `DataPipelineDefaultResourceRole`.

Beispiel für eine Richtlinie über Berechtigungen für Pipeline-Rollen

Die folgende Beispielrichtlinie ist darauf ausgelegt, wichtige Funktionen zuzulassen, die den Betrieb einer Pipeline mit Amazon EC2 - und Amazon EMR-Ressourcen AWS Data Pipeline erfordern. Es bietet auch Berechtigungen für den Zugriff auf andere AWS Ressourcen wie Amazon Simple Storage Service und Amazon Simple Notification Service, die für viele Pipelines erforderlich sind. Wenn die in einer Pipeline definierten Objekte nicht die Ressourcen eines AWS Dienstes benötigen, empfehlen wir dringend, dass Sie die Zugriffsberechtigungen für diesen Service entfernen. Wenn Ihre Pipeline beispielsweise keine Aktion definiert [Dynamo-Knoten DBData](#) oder die [SnsAlarm](#) Aktion nicht verwendet, empfehlen wir Ihnen, die Allow-Anweisungen für diese Aktionen zu entfernen.

- **111122223333** Ersetzen Sie es durch Ihre AWS Konto-ID.
- **NameOfDataPipelineRole** Ersetzen Sie es durch den Namen der Pipeline-Rolle (die Rolle, der diese Richtlinie zugeordnet ist).
- Durch **NameOfDataPipelineResourceRole** den Namen der EC2 Instanzrolle ersetzen.
- **us-west-1** Ersetzen Sie durch die entsprechende Region für Ihre Anwendung.

Verwaltete Standardrichtlinie für die EC2 Instanzrolle

Der Inhalt von `AmazonEC2RoleforDataPipelineRole` ist unten dargestellt. Dies ist die verwaltete Richtlinie, die der Standardressourcenrolle für AWS Data Pipeline, zugeordnet ist `AmazonDataPipelineDefaultResourceRole`. Wenn Sie eine Ressourcenrolle für Ihre Pipeline definieren, empfehlen wir, mit dieser Berechtigungsrichtlinie zu beginnen und dann Berechtigungen für AWS Dienstaktionen zu entfernen, die nicht erforderlich sind.

Version 3 der Richtlinie wird angezeigt. Dies ist die neueste Version zum Zeitpunkt der Erstellung dieses Artikels. Zeigen Sie die neueste Version der Richtlinie mithilfe der IAM-Konsole an.

JSON

```
{  
    "Version": "2012-10-17",  
    "Statement": [{  
        "Effect": "Allow",  
        "Action": [  
            "cloudwatch:*",  
            "datapipeline:*",  
            "dynamodb:*",  
            "ec2:Describe*",  
            "elasticmapreduce:AddJobFlowSteps",  
            "elasticmapreduce:Describe*",  
            "elasticmapreduce>ListInstance*",  
            "elasticmapreduce:ModifyInstanceGroups",  
            "rds:Describe*",  
            "redshift:DescribeClusters",  
            "redshift:DescribeClusterSecurityGroups",  
            "s3:*",  
            "sdb:*",  
            "sns:*",  
            "sqS:*"  
        ],  
        "Resource": ["*"]  
    }]  
}
```

IAM-Rollen für Rollenberechtigungen erstellen AWS Data Pipeline und bearbeiten

Gehen Sie wie folgt vor, um Rollen für die AWS Data Pipeline Verwendung der IAM-Konsole zu erstellen. Der Prozess besteht aus zwei Schritten. Zunächst erstellen Sie eine Berechtigungsrichtlinie, die Sie der Rolle zuordnen möchten. Als Nächstes erstellen Sie die Rolle und fügen die Richtlinie hinzu. Nachdem Sie eine Rolle erstellt haben, können Sie die Berechtigungen der Rolle ändern, indem Sie Berechtigungsrichtlinien anhängen oder trennen.

Note

Wenn Sie wie unten beschrieben Rollen für die AWS Data Pipeline Verwendung der Konsole erstellen, erstellt IAM die entsprechenden Vertrauensrichtlinien, die für die Rolle erforderlich sind, und fügt sie an.

Um eine Berechtigungsrichtlinie zur Verwendung mit einer Rolle für zu erstellen AWS Data Pipeline

1. Öffnen Sie unter <https://console.aws.amazon.com/iam/> die IAM-Konsole.
2. Wählen Sie im Navigationsbereich Richtlinien und dann Richtlinie erstellen.
3. Wählen Sie den Tab JSON.
4. Wenn Sie eine Pipeline-Rolle erstellen, kopieren Sie den Inhalt des Richtlinienbeispiels [Beispiel für eine Richtlinie über Berechtigungen für Pipeline-Rollen](#), fügen Sie ihn ein und bearbeiten Sie ihn entsprechend Ihren Sicherheitsanforderungen. Wenn Sie eine benutzerdefinierte EC2 Instanzrolle erstellen, tun Sie alternativ dasselbe für das Beispiel unter [Verwaltete Standardrichtlinie für die EC2 Instanzrolle](#).
5. Wählen Sie Richtlinie prüfen.
6. Geben Sie beispielsweise einen Namen für die Richtlinie und optional eine MyDataPipelineRolePolicy Beschreibung ein und wählen Sie dann Richtlinie erstellen aus.
7. Notieren Sie sich den Namen der Richtlinie. Sie benötigen ihn, wenn Sie Ihre Rolle erstellen.

Um eine IAM-Rolle zu erstellen für AWS Data Pipeline

1. Öffnen Sie unter <https://console.aws.amazon.com/iam/> die IAM-Konsole.
2. Wählen Sie im Navigationsbereich Rollen und dann Rolle erstellen aus.
3. Wählen Sie unter Anwendungsfall auswählen die Option Data Pipeline aus.

4. Führen Sie unter Wählen Sie Ihren Anwendungsfall aus eine der folgenden Aktionen aus:
 - Wählen Sie Data Pipeline, ob Sie eine Pipeline-Rolle erstellen möchten.
 - Wählen Sie EC2 Role for Data Pipeline, ob Sie eine Ressourcenrolle erstellen möchten.
5. Wählen Sie Weiter: Berechtigungen aus.
6. Wenn die Standardrichtlinie für aufgeführt AWS Data Pipeline ist, fahren Sie mit den folgenden Schritten fort, um die Rolle zu erstellen, und bearbeiten Sie sie dann gemäß den Anweisungen im nächsten Verfahren. Geben Sie andernfalls den Namen der Richtlinie ein, die Sie im obigen Verfahren erstellt haben, und wählen Sie sie dann aus der Liste aus.
7. Wählen Sie Weiter: Stichwörter aus, geben Sie alle Tags ein, die Sie der Rolle hinzufügen möchten, und wählen Sie dann Weiter: Überprüfen aus.
8. Geben Sie beispielsweise einen Namen für die Rolle und optional eine Beschreibung **MyDataPipelineRole** ein und wählen Sie dann Rolle erstellen aus.

Um eine Berechtigungsrichtlinie für eine IAM-Rolle anzuhängen oder zu trennen für AWS Data Pipeline

1. Öffnen Sie unter <https://console.aws.amazon.com/iam/> die IAM-Konsole.
2. Wählen Sie im Navigationsbereich Rollen
3. Geben Sie im Suchfeld den Namen der Rolle ein, die Sie bearbeiten möchten, z. B. DataPipelineDefaultRole oder, und wählen Sie MyDataPipelineRole dann den Rollennamen aus der Liste aus.
4. Gehen Sie auf der Registerkarte Berechtigungen wie folgt vor:
 - Um eine Berechtigungsrichtlinie zu trennen, klicken Sie unter Berechtigungsrichtlinien auf die Schaltfläche Entfernen ganz rechts neben dem Richtlinieneintrag. Wählen Sie Trennen, wenn Sie zur Bestätigung aufgefordert werden.
 - Um eine Richtlinie anzuhängen, die Sie zuvor erstellt haben, wählen Sie Richtlinien anhängen. Geben Sie im Suchfeld den Namen der Richtlinie ein, die Sie bearbeiten möchten, wählen Sie sie aus der Liste aus und wählen Sie dann Richtlinie anhängen aus.

Rollen für eine bestehende Pipelines ändern

Wenn Sie einer Pipeline eine andere Pipeline-Rolle oder Ressourcenrolle zuweisen möchten, können Sie den Architekt-Editor in der AWS Data Pipeline Konsole verwenden.

Um die einer Pipeline zugewiesenen Rollen mithilfe der Konsole zu bearbeiten

1. Öffnen Sie die AWS Data Pipeline Konsole unter <https://console.aws.amazon.com/datapipeline/>.
2. Wählen Sie die Pipeline aus der Liste aus und klicken Sie dann auf Aktionen, Bearbeiten.
3. Wählen Sie im rechten Bereich des Architect-Editors die Option Andere aus.
4. Wählen Sie aus den Listen „Ressourcenrolle“ und „Rolle“ die Rollen aus AWS Data Pipeline , denen Sie zuweisen möchten, und klicken Sie dann auf Speichern.

Einloggen und Überwachen AWS Data Pipeline

AWS Data Pipeline ist in einen Dienst integriert AWS CloudTrail, der eine Aufzeichnung der Aktionen bereitstellt, die von einem Benutzer, einer Rolle oder einem AWS Dienst in ausgeführt wurden AWS Data Pipeline. CloudTrail erfasst alle API-Aufrufe AWS Data Pipeline als Ereignisse. Zu den erfassten Aufrufen gehören Aufrufe von der AWS Data Pipeline Konsole und Codeaufrufen für die AWS Data Pipeline API-Operationen. Wenn Sie einen Trail erstellen, können Sie die kontinuierliche Bereitstellung von CloudTrail Ereignissen an einen Amazon S3 S3-Bucket aktivieren, einschließlich Ereignissen für AWS Data Pipeline. Wenn Sie keinen Trail konfigurieren, können Sie die neuesten Ereignisse trotzdem in der CloudTrail Konsole im Ereignisverlauf anzeigen. Anhand der von gesammelten Informationen können Sie die Anfrage ermitteln CloudTrail, an die die Anfrage gestellt wurde AWS Data Pipeline, die IP-Adresse, von der aus die Anfrage gestellt wurde, wer die Anfrage gestellt hat, wann sie gestellt wurde, und weitere Details.

Weitere Informationen CloudTrail dazu finden Sie im [AWS CloudTrail Benutzerhandbuch](#).

AWS Data Pipeline Informationen in CloudTrail

CloudTrail ist in Ihrem AWS Konto aktiviert, wenn Sie das Konto erstellen. Wenn eine Aktivität in stattfindet AWS Data Pipeline, wird diese Aktivität zusammen mit anderen CloudTrail AWS Serviceereignissen im Ereignisverlauf in einem Ereignis aufgezeichnet. Sie können die neusten Ereignisse in Ihr AWS -Konto herunterladen und dort suchen und anzeigen. Weitere Informationen finden Sie unter [Ereignisse mit CloudTrail Ereignisverlauf anzeigen](#).

Für eine fortlaufende Aufzeichnung der Ereignisse in Ihrem AWS Konto, einschließlich der Ereignisse für AWS Data Pipeline, erstellen Sie einen Trail. Ein Trail ermöglicht CloudTrail die Übermittlung von Protokolldateien an einen Amazon S3 S3-Bucket. Wenn Sie einen Pfad in der Konsole anlegen, gilt dieser für alle AWS-Regionen. Der Trail protokolliert Ereignisse aus allen Regionen der AWS Partition und übermittelt die Protokolldateien an den von Ihnen angegebenen Amazon S3 S3-Bucket. Darüber hinaus können Sie andere AWS Dienste konfigurieren, um die in den CloudTrail Protokollen gesammelten Ereignisdaten weiter zu analysieren und darauf zu reagieren. Weitere Informationen finden Sie hier:

- [Übersicht zum Erstellen eines Trails](#)
- [CloudTrail Unterstützte Dienste und Integrationen](#)
- [Konfiguration von Amazon SNS SNS-Benachrichtigungen für CloudTrail](#)
- [Empfangen von CloudTrail Protokolldateien aus mehreren Regionen](#) und [Empfangen von CloudTrail Protokolldateien von mehreren Konten](#)

Alle AWS Data Pipeline Aktionen werden im Kapitel AWS Data Pipeline API Reference Actions protokolliert CloudTrail und sind im [Kapitel AWS Data Pipeline API Reference Actions](#) dokumentiert. Beispielsweise generieren Aufrufe der CreatePipelineAktion Einträge in den CloudTrail Protokolldateien.

Jeder Ereignis- oder Protokolleintrag enthält Informationen zu dem Benutzer, der die Anforderung generiert hat. Die Identitätsinformationen unterstützen Sie bei der Ermittlung der folgenden Punkte:

- Ob die Anfrage mit Root- oder IAM-Rollenanmeldedaten gestellt wurde.
- Gibt an, ob die Anforderung mit temporären Sicherheitsanmeldeinformationen für eine Rolle oder einen Verbundbenutzer gesendet wurde.
- Ob die Anfrage von einem anderen AWS Dienst gestellt wurde.

Weitere Informationen finden Sie unter [CloudTrail userIdentity-Element](#).

Grundlegendes zu AWS Data Pipeline Protokolldateieinträgen

Ein Trail ist eine Konfiguration, die die Übertragung von Ereignissen als Protokolldateien an einen von Ihnen angegebenen Amazon S3 S3-Bucket ermöglicht. CloudTrail Protokolldateien enthalten einen oder mehrere Protokolleinträge. Ein Ereignis stellt eine einzelne Anforderung aus einer beliebigen Quelle dar und enthält Informationen über die angeforderte Aktion, Datum und Uhrzeit der Aktion,

Anforderungsparameter usw. CloudTrail Protokolldateien sind kein geordneter Stack-Trace der öffentlichen API-Aufrufe, sodass sie nicht in einer bestimmten Reihenfolge angezeigt werden.

Das folgende Beispiel zeigt einen CloudTrail Protokolleintrag, der den CreatePipeline Vorgang demonstriert:

```
{  
  "Records": [  
    {  
      "eventVersion": "1.02",  
      "userIdentity": {  
        "type": "Root",  
        "principalId": "123456789012",  
        "arn": "arn:aws:iam::aws-account-id:role/role-name",  
        "accountId": "role-account-id",  
        "accessKeyId": "role-access-key"  
      },  
      "eventTime": "2014-11-13T19:15:15Z",  
      "eventSource": "datapipeline.amazonaws.com",  
      "eventName": "CreatePipeline",  
      "awsRegion": "us-east-1",  
      "sourceIPAddress": "72.21.196.64",  
      "userAgent": "aws-cli/1.5.2 Python/2.7.5 Darwin/13.4.0",  
      "requestParameters": {  
        "name": "testpipeline",  
        "uniqueId": "sounique"  
      },  
      "responseElements": {  
        "pipelineId": "df-06372391ZG65EXAMPLE"  
      },  
      "requestID": "65cbf1e8-6b69-11e4-8816-cfcbad04c45",  
      "eventID": "9f99dce0-0864-49a0-bffa-f72287197758",  
      "eventType": "AwsApiCall",  
      "recipientAccountId": "role-account-id"  
    },  
    ...additional entries  
  ]  
}
```

Reaktion auf Vorfälle in AWS Data Pipeline

Die Reaktion auf Vorfälle für AWS Data Pipeline ist eine AWS Verantwortung. AWS verfügt über eine formelle, dokumentierte Richtlinie und ein Programm, die die Reaktion auf Vorfälle regeln.

Operative AWS-Probleme mit weitreichenden Auswirkungen werden im AWS Service Health Dashboard veröffentlicht. Operative Probleme werden ebenfalls über das Personal Health Dashboard in den einzelnen Konten gepostet.

Überprüfung der Einhaltung von Vorschriften für AWS Data Pipeline

AWS Data Pipeline fällt nicht in den Geltungsbereich von AWS-Compliance-Programmen. Eine Liste der AWS-Services, die in bestimmten Compliance-Programmen enthalten sind, finden Sie unter [AWS Services in Scope nach Compliance-Programm](#). Allgemeine Informationen finden Sie unter [AWS-Compliance-Programme](#).

Resilienz in AWS Data Pipeline

Die AWS globale Infrastruktur basiert auf AWS Regionen und Availability Zones. AWS Regionen bieten mehrere physisch getrennte und isolierte Availability Zones, die über Netzwerke mit niedriger Latenz, hohem Durchsatz und hoher Redundanz miteinander verbunden sind. Mithilfe von Availability Zones können Sie Anwendungen und Datenbanken erstellen und ausführen, die automatisch Failover zwischen Zonen ausführen, ohne dass es zu Unterbrechungen kommt. Availability Zones sind besser verfügbar, fehlertoleranter und skalierbarer als herkömmliche Infrastrukturen mit einem oder mehreren Rechenzentren.

Weitere Informationen zu AWS Regionen und Availability Zones finden Sie unter [AWS Globale Infrastruktur](#).

Sicherheit der Infrastruktur in AWS Data Pipeline

Als verwalteter Service AWS Data Pipeline ist er durch die AWS globalen Netzwerksicherheitsverfahren geschützt, die im Whitepaper [Amazon Web Services: Sicherheitsprozesse im Überblick](#) beschrieben sind.

Sie verwenden AWS veröffentlichte API-Aufrufe für den Zugriff AWS Data Pipeline über das Netzwerk. Kunden müssen Transport Layer Security (TLS) 1.0 oder neuer unterstützen. Wir

empfehlen TLS 1.2 oder höher. Clients müssen außerdem Verschlüsselungssammlungen mit PFS (Perfect Forward Secrecy) wie DHE (Ephemeral Diffie-Hellman) oder ECDHE (Elliptic Curve Ephemeral Diffie-Hellman) unterstützen. Die meisten modernen Systemen wie Java 7 und höher unterstützen diese Modi.

Außerdem müssen Anforderungen mit einer Zugriffsschlüssel-ID und einem geheimen Zugriffsschlüssel signiert sein, der einem IAM-Prinzipal zugeordnet ist. Alternativ können Sie mit [AWS Security Token Service](#) (AWS STS) temporäre Sicherheitsanmeldeinformationen erstellen, um die Anforderungen zu signieren.

Konfiguration und Schwachstellenanalyse in AWS Data Pipeline

Konfiguration und IT-Kontrollen liegen in der gemeinsamen Verantwortung AWS von Ihnen, unserem Kunden. Weitere Informationen finden Sie im [Modell der AWS gemeinsamen Verantwortung](#).

Tutorials

Die folgenden Tutorials führen Sie step-by-step durch den Prozess der Erstellung und Verwendung von Pipelines mit AWS Data Pipeline

Tutorials

- [Daten mithilfe von Amazon EMR mit Hadoop Streaming verarbeiten](#)
- [Kopieren Sie CSV-Daten zwischen Amazon S3 S3-Buckets mithilfe AWS Data Pipeline](#)
- [Exportieren von MySQL-Daten nach Amazon S3 mithilfe AWS Data Pipeline](#)
- [Daten mithilfe von Amazon Redshift nach Amazon Redshift kopieren AWS Data Pipeline](#)

Daten mithilfe von Amazon EMR mit Hadoop Streaming verarbeiten

Sie können es AWS Data Pipeline zur Verwaltung Ihrer Amazon EMR-Cluster verwenden. Mit können AWS Data Pipeline Sie Vorbedingungen angeben, die erfüllt sein müssen, bevor der Cluster gestartet wird (z. B. sicherstellen, dass die heutigen Daten auf Amazon S3 hochgeladen wurden), einen Zeitplan für die wiederholte Ausführung des Clusters und die zu verwendende Cluster-Konfiguration angeben. Das folgende Tutorial führt Sie durch den Start eines einfachen Clusters.

In diesem Tutorial erstellen Sie eine Pipeline für einen einfachen Amazon EMR-Cluster, um einen bereits vorhandenen Hadoop-Streaming-Job auszuführen, der von Amazon EMR bereitgestellt wird, und eine Amazon SNS SNS-Benachrichtigung zu senden, nachdem die Aufgabe erfolgreich abgeschlossen wurde. Für diese Aufgabe verwenden Sie die Amazon EMR-Cluster-Ressource, AWS Data Pipeline die von bereitgestellt wird. Die Beispielanwendung wird aufgerufen WordCount und kann auch manuell von der Amazon EMR-Konsole aus ausgeführt werden. Beachten Sie, dass Cluster, die von in AWS Data Pipeline Ihrem Namen erzeugt wurden, in der Amazon EMR-Konsole angezeigt und Ihrem AWS-Konto in Rechnung gestellt werden.

Pipeline-Objekte

Die Pipeline verwendet die folgenden Objekte:

[EmrActivity](#)

Definiert die Arbeit, die in der Pipeline ausgeführt werden soll (einen bereits vorhandenen Hadoop-Streaming-Job ausführen, der von Amazon EMR bereitgestellt wird).

EmrCluster

Die Ressource, AWS Data Pipeline die zur Ausführung dieser Aktivität verwendet wird.

Ein Cluster besteht aus einer Reihe von EC2 Amazon-Instances. AWS Data Pipeline startet den Cluster und beendet ihn dann, nachdem die Aufgabe abgeschlossen ist.

Plan

Startdatum, Uhrzeit und Dauer dieser Aktivität. Sie können optional das Enddatum und die Endzeit angeben.

SnsAlarm

Sendet eine Amazon SNS SNS-Benachrichtigung an das von Ihnen angegebene Thema, nachdem die Aufgabe erfolgreich abgeschlossen wurde.

Inhalt

- [Bevor Sie beginnen](#)
- [Einen Cluster über die Befehlszeile starten](#)

Bevor Sie beginnen

Stellen Sie sicher, dass Sie die folgenden Schritte ausgeführt haben.

- Führen Sie die Aufgaben unter [Einrichtung für AWS Data Pipeline](#).
- (Optional) Richten Sie eine VPC für den Cluster und eine Sicherheitsgruppe für die VPC ein.
- Erstellen Sie ein Thema zum Senden einer E-Mail-Benachrichtigung und notieren Sie sich das Thema Amazon-Ressourcename (ARN). Weitere Informationen dazu erhalten Sie unter [Erstellen eines Themas](#) im Amazon Simple Notification Service Handbuch Erste Schritte.

Einen Cluster über die Befehlszeile starten

Wenn Sie regelmäßig einen Amazon EMR-Cluster zur Analyse von Webprotokollen oder zur Analyse wissenschaftlicher Daten ausführen, können Sie ihn AWS Data Pipeline zur Verwaltung Ihrer Amazon EMR-Cluster verwenden. Mit können Sie Vorbedingungen angeben AWS Data Pipeline, die erfüllt sein müssen, bevor der Cluster gestartet wird (z. B. um sicherzustellen, dass die heutigen Daten auf

Amazon S3 hochgeladen wurden). Dieses Tutorial führt Sie durch die Einführung eines Clusters, der als Modell für eine einfache Amazon EMR-basierte Pipeline oder als Teil einer komplexeren Pipeline dienen kann.

Voraussetzungen

Bevor Sie die Befehlszeile zum ersten Mal verwenden können, müssen Sie die folgenden Schritte ausführen:

1. Installieren und konfigurieren Sie eine Befehlszeilenschnittstelle (CLI). Weitere Informationen finden Sie unter [Zugreifen AWS Data Pipeline](#).
2. Stellen Sie sicher, dass die IAM-Rollen benannt DataPipelineDefaultResourceRolesind DataPipelineDefaultRoleund existieren. Die AWS Data Pipeline Konsole erstellt diese Rollen automatisch für Sie. Wenn Sie die AWS Data Pipeline Konsole nicht mindestens einmal verwendet haben, müssen Sie diese Rollen manuell erstellen. Weitere Informationen finden Sie unter [IAM-Rollen für AWS Data Pipeline](#).

Aufgaben

- [Erstellen der Pipeline-Definitionsdatei](#)
- [Hochladen und Aktivieren der Pipeline-Definition](#)
- [Überwachen der Pipeline-Runs](#)

Erstellen der Pipeline-Definitionsdatei

Der folgende Code ist die Pipeline-Definitionsdatei für einen einfachen Amazon EMR-Cluster, der einen vorhandenen Hadoop-Streaming-Job ausführt, der von Amazon EMR bereitgestellt wird. Diese Beispielanwendung wird aufgerufen WordCount, und Sie können sie auch über die Amazon EMR-Konsole ausführen.

Kopieren Sie diesen Code in eine Textdatei, und speichern Sie sie unter `MyEmrPipelineDefinition.json`. Sie sollten den Amazon S3 S3-Bucket-Speicherort durch den Namen eines Amazon S3 S3-Buckets ersetzen, den Sie besitzen. Sie sollten auch das Start- und das Enddatum ersetzen. Um Cluster sofort zu starten, legen Sie `startDate` ein Datum fest, das einen Tag in der Vergangenheit und `endDate` einen Tag in der future hat. AWS Data Pipeline beginnt dann sofort mit dem Start der „überfälligen“ Cluster, um den als Arbeitsstau wahrgenommenen Arbeitsstau zu beheben. Dieses Backfilling bedeutet, dass Sie nicht eine Stunde warten müssen, bis der erste Cluster AWS Data Pipeline gestartet wird.

```
{  
  "objects": [  
    {  
      "id": "Hourly",  
      "type": "Schedule",  
      "startDateTime": "2012-11-19T07:48:00",  
      "endDateTime": "2012-11-21T07:48:00",  
      "period": "1 hours"  
    },  
    {  
      "id": "MyCluster",  
      "type": "EmrCluster",  
      "masterInstanceType": "m1.small",  
      "schedule": {  
        "ref": "Hourly"  
      }  
    },  
    {  
      "id": "MyEmrActivity",  
      "type": "EmrActivity",  
      "schedule": {  
        "ref": "Hourly"  
      },  
      "runsOn": {  
        "ref": "MyCluster"  
      },  
      "step": "/home/hadoop/contrib/streaming/hadoop-streaming.jar,-input,s3n://  
elasticmapreduce/samples/wordcount/input,-output,s3://myawsbucket/wordcount/  
output/#{@scheduledStartTime},-mapper,s3n://elasticmapreduce/samples/wordcount/  
wordSplitter.py,-reducer,aggregate"  
    }  
  ]  
}
```

Diese Pipeline hat drei Objekte:

- Hourly, was den Zeitplan für die Arbeit repräsentiert. Sie können einen Zeitplan als eines der Felder für eine Aktivität festlegen. Wenn Sie das tun, wird die Aktivität gemäß diesem Zeitplan ausgeführt, in diesem Fall stündlich.
- MyCluster, was die Gruppe von EC2 Amazon-Instances darstellt, die für den Betrieb des Clusters verwendet werden. Sie können die Größe und Anzahl der EC2 Instances angeben, die als Cluster ausgeführt werden sollen. Wenn Sie die Anzahl an Instances nicht festlegen, startet

der Cluster mit zwei Instances, einem Master-Knoten und einem Aufgabenknoten. Sie können ein Subnetz angeben, in dem der Cluster gestartet werden soll. Sie können dem Cluster zusätzliche Konfigurationen hinzufügen, z. B. Bootstrap-Aktionen, um zusätzliche Software auf das von Amazon EMR bereitgestellte AMI zu laden.

- MyEmrActivity, was die Berechnung darstellt, die mit dem Cluster verarbeitet werden soll. Amazon EMR unterstützt verschiedene Clustertypen, darunter Streaming, Cascading und Scripted Hive. Das runsOn Feld bezieht sich auf MyCluster und verwendet dies als Spezifikation für die Grundlagen des Clusters.

Hochladen und Aktivieren der Pipeline-Definition

Sie müssen Ihre Pipeline-Definition hochladen und Ihre Pipeline aktivieren. Ersetzen Sie die Befehle in den folgenden Beispielbefehlen *pipeline_name* durch ein Label für Ihre Pipeline und *pipeline_file* durch den vollqualifizierten Pfad für die .json Pipeline-Definitionsdatei.

AWS CLI

Verwenden Sie den folgenden Befehl [create-pipeline](#), um Ihre Pipeline-Definition zu erstellen und Ihre Pipeline zu aktivieren. Notieren Sie sich die ID Ihrer Pipeline, da Sie diesen Wert mit den meisten CLI-Befehlen verwenden werden.

```
aws datapipeline create-pipeline --name pipeline_name --unique-id token
{
    "pipelineId": "df-00627471S0VYZEXAMPLE"
}
```

Verwenden Sie den folgenden [put-pipeline-definition](#)Befehl, um Ihre Pipeline-Definition hochzuladen.

```
aws datapipeline put-pipeline-definition --pipeline-id df-00627471S0VYZEXAMPLE --
pipeline-definition file://MyEmrPipelineDefinition.json
```

Wenn Ihre Pipeline erfolgreich validiert wurde, ist das validationErrors Feld leer. Sie sollten alle Warnungen überprüfen.

Verwenden Sie den folgenden Befehl [activate-pipeline](#), um Ihre Pipeline zu aktivieren.

```
aws datapipeline activate-pipeline --pipeline-id df-00627471S0VYZEXAMPLE
```

Mit dem folgenden Befehl [list-pipelines](#) können Sie überprüfen, ob Ihre Pipeline in der Pipeline-Liste aufgeführt ist.

```
aws datapipeline list-pipelines
```

Überwachen der Pipeline-Runs

Sie können Cluster anzeigen, die AWS Data Pipeline mit der Amazon EMR-Konsole gestartet wurden, und Sie können den Ausgabeordner mit der Amazon S3 S3-Konsole anzeigen.

Um den Fortschritt von Clustern zu überprüfen, die gestartet wurden von AWS Data Pipeline

1. Öffnen Sie die Amazon-EMR-Konsole.
2. Die Cluster, die von erzeugt wurden, AWS Data Pipeline haben einen Namen, der wie folgt formatiert ist: <pipeline-identifier> @_<emr-cluster-name>. <launch-time>

The screenshot shows the 'Cluster List' section of the AWS EMR console. At the top, there are buttons for 'Create cluster', 'View details', 'Clone', and 'Terminate'. Below that is a filter dropdown set to 'All clusters' and a search bar. It displays '3 clusters (all loaded)'. A table lists two clusters:

Name	ID	Status
df-00592868ZT33HX1F5I0_@MyCluster_2014-06-29T02:00:00	j-20XJRAX6Z5HC4	Running
df-00592868ZT33HX1F5I0_@MyCluster_2014-06-29T01:00:00	j-32CYSLG57E6YT	Running

3. Nachdem einer der Läufe abgeschlossen ist, öffnen Sie die Amazon S3 S3-Konsole und überprüfen Sie, ob der Ausgabeordner mit Zeitstempel vorhanden ist und die erwarteten Ergebnisse des Clusters enthält.

The screenshot shows the 'wordcount' folder within the 'js-s3-bucket'. At the top, there are buttons for 'Upload', 'Create Folder', and 'Actions'. The path 'All Buckets / js-s3-bucket / wordcount' is shown. A table lists three objects:

Name
2014-06-29T00:00:00
2014-06-29T01:00:00
2014-06-29T02:00:00

Kopieren Sie CSV-Daten zwischen Amazon S3 S3-Buckets mithilfe AWS Data Pipeline

Nachdem Sie das gelesen haben [Was ist AWS Data Pipeline?](#) und entschieden haben, dass Sie die Übertragung und Transformation Ihrer Daten automatisieren AWS Data Pipeline möchten, ist es an der Zeit, mit der Erstellung von Daten-Pipelines zu beginnen. Um Ihnen die Funktionsweise von AWS Data Pipeline näher zu bringen, hier ein Beispiel für eine einfache Aufgabe.

Dieses Tutorial führt Sie durch den Prozess der Erstellung einer Datenpipeline, um Daten von einem Amazon S3 S3-Bucket in einen anderen zu kopieren und dann eine Amazon SNS SNS-Benachrichtigung zu senden, nachdem der Kopiervorgang erfolgreich abgeschlossen wurde. Für diese Kopieraktivität verwenden Sie eine EC2 Instance, AWS Data Pipeline die von verwaltet wird.

Pipeline-Objekte

Die Pipeline verwendet die folgenden Objekte:

[CopyActivity](#)

Die Aktivität, die AWS Data Pipeline für diese Pipeline ausgeführt wird (CSV-Daten von einem Amazon-S3-Bucket in einen anderen kopieren).

Important

Es gibt Einschränkungen bei der Verwendung des CSV-Dateiformats mit CopyActivity und S3DataNode. Weitere Informationen finden Sie unter [CopyActivity](#).

[Plan](#)

Das Startdatum, die Uhrzeit und die Wiederholung für diese Aktivität. Sie können optional das Enddatum und die Endzeit angeben.

[Ec2Resource](#)

Die Ressource (eine EC2 Instance), die zur Ausführung dieser Aktivität AWS Data Pipeline verwendet wird.

[S3 DataNode](#)

Die Eingabe- und Ausgabeknoten (Amazon S3 S3-Buckets) für diese Pipeline.

SnsAlarm

Die Aktion AWS Data Pipeline muss ausgeführt werden, wenn die angegebenen Bedingungen erfüllt sind (Amazon SNS SNS-Benachrichtigungen an ein Thema senden, nachdem die Aufgabe erfolgreich abgeschlossen wurde).

Inhalt

- [Bevor Sie beginnen](#)
- [CSV-Daten mithilfe der Befehlszeile kopieren](#)

Bevor Sie beginnen

Stellen Sie sicher, dass Sie die folgenden Schritte ausgeführt haben.

- Führen Sie die Aufgaben unter au [Einrichtung für AWS Data Pipeline](#).
- (Optional) Richten Sie eine VPC für die Instance und eine Sicherheitsgruppe für die VPC ein.
- Erstellen Sie einen Amazon-S3-Bucket als Datenquelle.

Weitere Informationen finden Sie unter [Erstellen eines Buckets](#) im Benutzerhandbuch zu Amazon Simple Storage Service.

- Laden Sie Ihre Daten in Ihren Amazon-S3-Bucket hoch.

Anleitungen finden Sie unter [Hinzufügen eines Objekts zu einem Bucket](#) im Benutzerhandbuch zu Amazon Simple Storage Service.

- Einen weiteren Amazon-S3-Bucket als Datenziel erstellen
- Erstellen Sie ein Thema zum Senden einer E-Mail-Benachrichtigung und notieren Sie sich das Thema Amazon-Ressourcename (ARN). Weitere Informationen dazu erhalten Sie unter [Erstellen eines Themas](#) im Amazon Simple Notification Service Handbuch Erste Schritte.
- (Optional) Dieses Tutorial verwendet die Standard-IAM-Rollenrichtlinien, die von AWS Data Pipeline erstellt wurden. Wenn Sie lieber Ihre eigenen IAM-Rollenrichtlinien und Vertrauensbeziehungen erstellen und konfigurieren möchten, folgen Sie den Anweisungen unter[IAM-Rollen für AWS Data Pipeline](#).

CSV-Daten mithilfe der Befehlszeile kopieren

Sie können Pipelines erstellen und verwenden, um Daten von einem Amazon-S3-Bucket in einen anderen zu kopieren.

Voraussetzungen

Bevor Sie beginnen, müssen Sie die folgenden Schritte ausführen:

1. Installieren und konfigurieren Sie eine Befehlszeilenschnittstelle (Command Line Interface, CLI). Weitere Informationen finden Sie unter [Zugreifen AWS Data Pipeline](#).
2. Stellen Sie sicher, dass die IAM-Rollen benannt DataPipelineDefaultResourceRolesind DataPipelineDefaultRoleund existieren. Die AWS Data Pipeline Konsole erstellt diese Rollen automatisch für Sie. Wenn Sie die AWS Data Pipeline Konsole nicht mindestens einmal verwendet haben, müssen Sie diese Rollen manuell erstellen. Weitere Informationen finden Sie unter [IAM-Rollen für AWS Data Pipeline](#).

Aufgaben

- [Definieren Sie eine Pipeline im JSON-Format](#)
- [Hochladen und Aktivieren der Pipeline-Definition](#)

Definieren Sie eine Pipeline im JSON-Format

Dieses Beispielszenario zeigt, wie JSON-Pipeline-Definitionen und die AWS Data Pipeline CLI verwendet werden, um das Kopieren von Daten zwischen zwei Amazon S3 S3-Buckets in einem bestimmten Zeitintervall zu planen. Dies ist die vollständige Pipeline-Definition-JSON-Datei, gefolgt von einer Erläuterung für jeden ihrer Abschnitte.

Note

Wir empfehlen, dass Sie einen Texteditor verwenden, mit dem Sie die Syntax von JSON-formatierten Dateien überprüfen und die Datei mit der Dateierweiterung .json benennen können.

In diesem Beispiel überspringen wir aus Gründen der Übersichtlichkeit die optionalen Felder und zeigen nur erforderliche Felder an. Die vollständige Pipeline-JSON-Datei für dieses Beispiel lautet:

```
{  
  "objects": [  
    {  
      "id": "MySchedule",  
      "type": "Schedule",  
      "startDateTime": "2013-08-18T00:00:00",  
      "endDateTime": "2013-08-19T00:00:00",  
      "period": "1 day"  
    },  
    {  
      "id": "S3Input",  
      "type": "S3DataNode",  
      "schedule": {  
        "ref": "MySchedule"  
      },  
      "filePath": "s3://amzn-s3-demo-bucket/source/inputfile.csv"  
    },  
    {  
      "id": "S3Output",  
      "type": "S3DataNode",  
      "schedule": {  
        "ref": "MySchedule"  
      },  
      "filePath": "s3://amzn-s3-demo-bucket/destination/outputfile.csv"  
    },  
    {  
      "id": "MyEC2Resource",  
      "type": "Ec2Resource",  
      "schedule": {  
        "ref": "MySchedule"  
      },  
      "instanceType": "m1.medium",  
      "role": "DataPipelineDefaultRole",  
      "resourceRole": "DataPipelineDefaultResourceRole"  
    },  
    {  
      "id": "MyCopyActivity",  
      "type": "CopyActivity",  
      "runsOn": {  
        "ref": "MyEC2Resource"  
      },  
      "input": {  
        "ref": "S3Input"  
      }  
    }  
  ]  
}
```

```
        },
        "output": {
            "ref": "S3Output"
        },
        "schedule": {
            "ref": "MySchedule"
        }
    }
]
```

Plan

Die Pipeline definiert einen Zeitplan mit einem Start- und Enddatum sowie einem Zeitraum, um zu bestimmen, wie häufig die Aktivität in dieser Pipeline ausgeführt wird.

```
{
    "id": "MySchedule",
    "type": "Schedule",
    "startDateTime": "2013-08-18T00:00:00",
    "endDateTime": "2013-08-19T00:00:00",
    "period": "1 day"
},
```

Datenknoten

Als Nächstes definiert die DataNode S3-Eingabe-Pipeline-Komponente einen Speicherort für die Eingabedateien, in diesem Fall einen Amazon S3 S3-Bucket-Speicherort. Die DataNode S3-Eingabekomponente wird durch die folgenden Felder definiert:

```
{
    "id": "S3Input",
    "type": "S3DataNode",
    "schedule": {
        "ref": "MySchedule"
    },
    "filePath": "s3://example-bucket/source/inputfile.csv"
},
```

Id

Der benutzerdefinierte Name für den Eingabestandort (eine Bezeichnung nur für Ihre Referenz).

Typ

Der Pipeline-Komponententyp, der „S3DataNode“ ist, um dem Ort zu entsprechen, an dem sich die Daten befinden, in einem Amazon S3 S3-Bucket.

Plan

Ein Verweis auf die Zeitplankomponente, die wir in den vorherigen Zeilen der JSON-Datei mit der Bezeichnung „MySchedule“ erstellt haben.

Pfad

Der Pfad zu den Daten, die dem Datenknoten zugeordnet sind. Die Syntax für einen Datenknoten wird von seinem Typ bestimmt. Beispielsweise folgt die Syntax für einen Amazon S3 S3-Pfad einer anderen Syntax, die für eine Datenbanktabelle geeignet ist.

Als Nächstes definiert die DataNode Ausgabe-S3-Komponente den Ausgabzielort für die Daten. Sie folgt demselben Format wie die DataNode S3-Eingabekomponente, mit Ausnahme des Namens der Komponente und eines anderen Pfads zur Angabe der Zielfile.

```
{  
  "id": "S3Output",  
  "type": "S3DataNode",  
  "schedule": {  
    "ref": "MySchedule"  
  },  
  "filePath": "s3://example-bucket/destination/outputfile.csv"  
},
```

Ressource

Dies ist eine Definition der Rechenressource, die die Kopieroperation ausführt. In diesem Beispiel AWS Data Pipeline sollte automatisch eine EC2 Instanz zur Ausführung der Kopieraufgabe erstellt und die Ressource nach Abschluss der Aufgabe beendet werden. Die hier definierten Felder steuern die Erstellung und Funktion der EC2 Instanz, die die Arbeit erledigt. Die EC2 Ressource wird durch die folgenden Felder definiert:

```
{  
  "id": "MyEC2Resource",  
  "type": "Ec2Resource",  
  "schedule": {
```

```
    "ref": "MySchedule"
},
"instanceType": "m1.medium",
"role": "DataPipelineDefaultRole",
"resourceRole": "DataPipelineDefaultResourceRole"
},
```

Id

Der benutzerdefinierte Name für den Pipeline-Zeitplan, der nur für Ihre Referenz gilt.

Typ

Der Typ der Rechenressource, mit der Arbeit ausgeführt werden soll; in diesem Fall eine EC2 Instanz. Es sind auch andere Ressourcentypen verfügbar, z. B. ein EmrCluster Typ.

Plan

Der Zeitplan für die Erstellung dieser Rechenressource.

instanceType

Die Größe der zu erstellenden EC2 Instanz. Stellen Sie sicher, dass Sie die Größe der EC2 Instanz so einstellen, dass sie der Arbeitslast, mit der Sie arbeiten möchten, am besten entspricht AWS Data Pipeline. In diesem Fall haben wir eine EC2 m1.medium-Instanz eingerichtet. Weitere Informationen zu den verschiedenen Instance-Typen und zu den Einsatzmöglichkeiten der einzelnen Instance-Typen finden Sie im Thema [EC2Amazon-Instance-Typen](#) unter <http://aws.amazon.com/ec2/instance-types/>.

Rolle

Die IAM-Rolle des Kontos, das auf Ressourcen zugreift, z. B. auf einen Amazon S3 S3-Bucket zum Abrufen von Daten.

resourceRole

Die IAM-Rolle des Kontos, das Ressourcen erstellt, z. B. das Erstellen und Konfigurieren einer EC2 Instance in Ihrem Namen. Rolle und ResourceRole können dieselbe Rolle sein, bieten jedoch getrennt eine größere Granularität in Ihrer Sicherheitskonfiguration.

Aktivität

Der letzte Abschnitt in der JSON-Datei ist die Definition der Aktivität, die die auszuführende Arbeit darstellt. In diesem Beispiel CopyActivity werden Daten aus einer CSV-Datei in ein http://

aws.amazon kopiert. com/ec2/instance-types/ Bucket in einen anderen. Die CopyActivity-Komponente ist durch folgende Felder definiert:

```
{  
  "id": "MyCopyActivity",  
  "type": "CopyActivity",  
  "runsOn": {  
    "ref": "MyEC2Resource"  
  },  
  "input": {  
    "ref": "S3Input"  
  },  
  "output": {  
    "ref": "S3Output"  
  },  
  "schedule": {  
    "ref": "MySchedule"  
  }  
}
```

Id

Der benutzerdefinierte Name für die Aktivität, der nur für Ihre Referenz eine Bezeichnung ist.

Typ

Die Art der auszuführenden Aktivität, z. B. MyCopyActivity

runsOn

Die Datenverarbeitungsressource, die die Arbeit ausführt, die diese Aktivität definiert. In diesem Beispiel geben wir einen Verweis auf die zuvor definierte EC2 Instanz an. Wenn Sie das runsOn Feld verwenden AWS Data Pipeline , wird die EC2 Instanz für Sie erstellt. Das Feld runsOn zeigt an, dass die Ressource in der AWS-Infrastruktur vorhanden ist, während der Wert workerGroup angibt, dass Sie Ihre eigenen lokalen Ressourcen zur Ausführung der Arbeit verwenden möchten.

Eingabe

Der Speicherort der zu kopierenden Daten.

Output

Die Zielortdaten.

Plan

Der Zeitplan für die Ausführung dieser Aktivität.

Hochladen und Aktivieren der Pipeline-Definition

Sie müssen Ihre Pipeline-Definition hochladen und Ihre Pipeline aktivieren. Ersetzen Sie die Befehle in den folgenden Beispielbefehlen *pipeline_name* durch ein Label für Ihre Pipeline und *pipeline_file* durch den vollqualifizierten Pfad für die .json Pipeline-Definitionsdatei.

AWS CLI

Verwenden Sie den folgenden Befehl [create-pipeline](#), um Ihre Pipeline-Definition zu erstellen und Ihre Pipeline zu aktivieren. Notieren Sie sich die ID Ihrer Pipeline, da Sie diesen Wert mit den meisten CLI-Befehlen verwenden werden.

```
aws datapipeline create-pipeline --name pipeline_name --unique-id token
{
    "pipelineId": "df-00627471S0VYZEXAMPLE"
}
```

Verwenden Sie den folgenden [put-pipeline-definition](#)Befehl, um Ihre Pipeline-Definition hochzuladen.

```
aws datapipeline put-pipeline-definition --pipeline-id df-00627471S0VYZEXAMPLE --
pipeline-definition file://MyEmrPipelineDefinition.json
```

Wenn Ihre Pipeline erfolgreich validiert wurde, ist das validationErrors Feld leer. Sie sollten alle Warnungen überprüfen.

Verwenden Sie den folgenden Befehl [activate-pipeline](#), um Ihre Pipeline zu aktivieren.

```
aws datapipeline activate-pipeline --pipeline-id df-00627471S0VYZEXAMPLE
```

Mit dem folgenden Befehl [list-pipelines](#) können Sie überprüfen, ob Ihre Pipeline in der Pipeline-Liste aufgeführt ist.

```
aws datapipeline list-pipelines
```

Exportieren von MySQL-Daten nach Amazon S3 mithilfe AWS Data Pipeline

Dieses Tutorial führt Sie durch den Prozess der Erstellung einer Datenpipeline zum Kopieren von Daten (Zeilen) aus einer Tabelle in der MySQL-Datenbank in eine CSV-Datei (durch Kommas getrennte Werte) in einem Amazon S3-Bucket und zum Senden einer Amazon SNS-Benachrichtigung, nachdem die Kopieraktivität erfolgreich abgeschlossen wurde. Für diese Kopieraktivität verwenden Sie eine von AWS Data Pipeline bereitgestellte EC2 Instanz.

Pipeline-Objekte

Die Pipeline verwendet die folgenden Objekte:

- [CopyActivity](#)
- [Ec2Resource](#)
- [MySqlDataNode](#)
- [S3 DataNode](#)
- [SnsAlarm](#)

Inhalt

- [Bevor Sie beginnen](#)
- [MySQL-Daten über die Befehlszeile kopieren](#)

Bevor Sie beginnen

Stellen Sie sicher, dass Sie die folgenden Schritte ausgeführt haben.

- Führen Sie die Aufgaben unter [Einrichtung für AWS Data Pipeline](#).
- (Optional) Richten Sie eine VPC für die Instance und eine Sicherheitsgruppe für die VPC ein.
- Erstellen eines Amazon-S3-Buckets als Datenausgabe.

Weitere Informationen finden Sie unter [Create a Bucket](#) im Amazon Simple Storage Service-Benutzerhandbuch.

- Erstellen und starten Sie eine MySQL-Datenbank-Instance als Datenquelle.

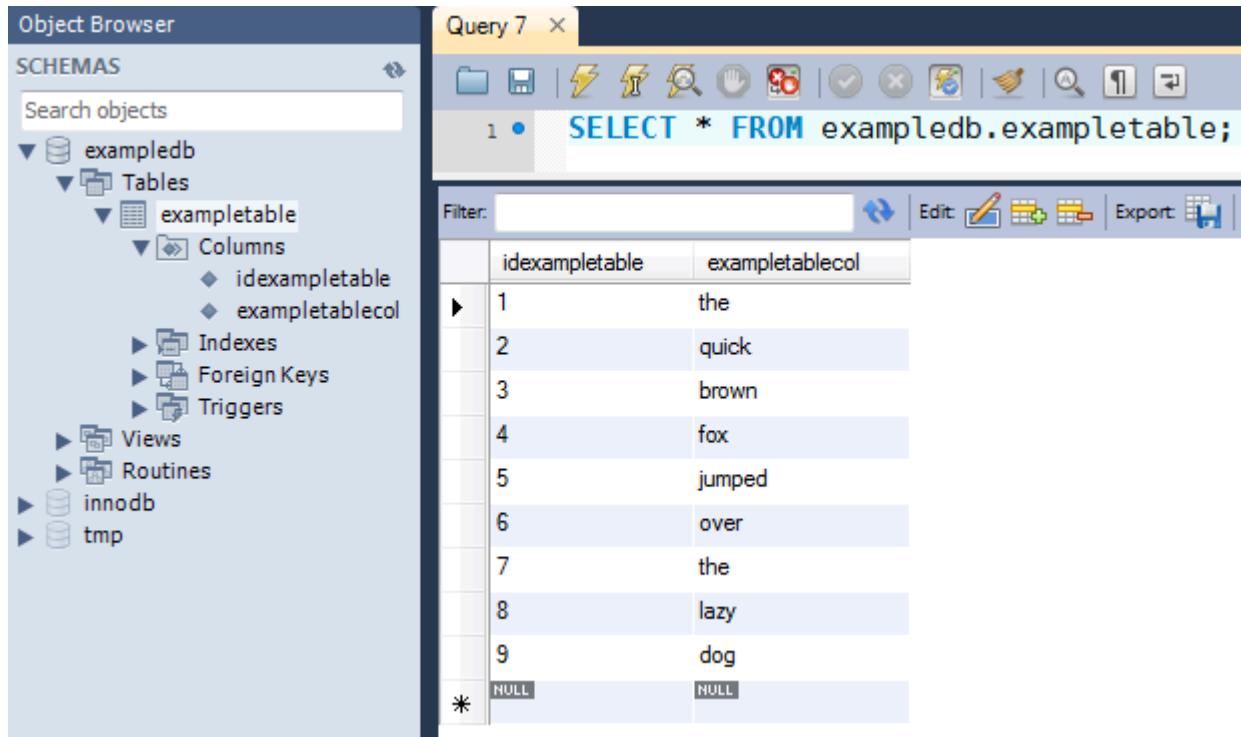
Weitere Informationen finden Sie unter [Launch a DB Instance](#) im Amazon RDS Getting Started Guide. Nachdem Sie eine Amazon RDS-Instance eingerichtet haben, finden Sie weitere Informationen unter [Create a Table](#) in der MySQL-Dokumentation.

 Note

Notieren Sie sich den Benutzernamen und das Passwort, das Sie beim Erstellen der MySQL-Instance verwendet haben. Nachdem Sie Ihre MySQL-Datenbank-Instance gestartet haben, notieren Sie sich den Endpunkt der Instance. Sie benötigen diese Informationen später wieder.

- Stellen Sie eine Verbindung mit der MySQL-Datenbank-Instance her, erstellen Sie eine Tabelle und fügen Sie die Testdatenwerte zur neu erstellten Tabelle hinzu.

Zur Verdeutlichung und Unterstützung haben wir dieses Tutorial erstellt, mit einer MySQL-Tabelle mit der folgenden Konfiguration und Beispieldaten. Der folgende Screenshot stammt von MySQL Workbench 5.2 CE:



The screenshot shows the MySQL Workbench interface. On the left, the Object Browser displays the database schema. It shows a 'SCHEMAS' section with 'exampledb' selected, containing 'Tables', 'Views', and 'Routines'. Under 'Tables', 'exampletable' is selected, showing its 'Columns' (idexamptable, examptablecol), 'Indexes', 'Foreign Keys', and 'Triggers'. On the right, the 'Query 7' tab is active, displaying the SQL query: 'SELECT * FROM exampledb.exampletable;'. The results are shown in a table with two columns: 'idexamptable' and 'examptablecol'. The data rows are: 1, the; 2, quick; 3, brown; 4, fox; 5, jumped; 6, over; 7, the; 8, lazy; 9, dog; followed by a row with NULL values.

	idexamptable	examptablecol
1		the
2		quick
3		brown
4		fox
5		jumped
6		over
7		the
8		lazy
9		dog
*		
	HULL	NULL

Weitere Informationen finden Sie unter [Erstellen einer Tabelle](#) in der MySQL-Dokumentation und auf der [MySQL Workbench-Produktseite](#).

- Erstellen Sie ein Thema zum Senden einer E-Mail-Benachrichtigung und notieren Sie sich das Thema Amazon-Ressourcename (ARN). Weitere Informationen finden Sie unter [Thema erstellen](#) im Amazon Simple Notification Service Getting Started Guide.
- (Optional) In diesem Tutorial werden die standardmäßigen IAM-Rollenrichtlinien verwendet, die von AWS Data Pipeline erstellt wurden. Wenn Sie lieber Ihre IAM-Rollenrichtlinie und Ihre Vertrauensbeziehungen erstellen und konfigurieren möchten, folgen Sie den Anweisungen unter [IAM-Rollen für AWS Data Pipeline](#)

MySQL-Daten über die Befehlszeile kopieren

Sie können eine Pipeline erstellen, um Daten aus einer MySQL-Tabelle in einen Amazon-S3-Bucket zu kopieren.

Voraussetzungen

Bevor Sie beginnen, müssen Sie die folgenden Schritte ausführen:

1. Installieren und konfigurieren Sie eine Befehlszeilenschnittstelle (Command Line Interface, CLI). Weitere Informationen finden Sie unter [Zugreifen AWS Data Pipeline](#).
2. Stellen Sie sicher, dass die IAM-Rollen benannt DataPipelineDefaultResourceRolesind DataPipelineDefaultRoleund existieren. Die AWS Data Pipeline Konsole erstellt diese Rollen automatisch für Sie. Wenn Sie die AWS Data Pipeline Konsole nicht mindestens einmal verwendet haben, müssen Sie diese Rollen manuell erstellen. Weitere Informationen finden Sie unter [IAM-Rollen für AWS Data Pipeline](#).
3. Richten Sie einen Amazon-S3-Bucket und eine Amazon-RDS-Instance ein. Weitere Informationen finden Sie unter [Bevor Sie beginnen](#).

Aufgaben

- [Definieren Sie eine Pipeline im JSON-Format](#)
- [Hochladen und Aktivieren der Pipeline-Definition](#)

Definieren Sie eine Pipeline im JSON-Format

Dieses Beispielszenario zeigt, wie JSON-Pipeline-Definitionen und die AWS Data Pipeline CLI verwendet werden, um Daten (Zeilen) aus einer Tabelle in einer MySQL-Datenbank in eine CSV-

Datei (durch Kommas getrennte Werte) in einem Amazon S3 S3-Bucket in einem bestimmten Zeitintervall zu kopieren.

Dies ist die vollständige Pipeline-Definition-JSON-Datei, gefolgt von einer Erläuterung für jeden ihrer Abschnitte.

 Note

Wir empfehlen, dass Sie einen Texteditor verwenden, mit dem Sie die Syntax von JSON-formatierten Dateien überprüfen und die Datei mit der Dateierweiterung .json benennen können.

```
{  
  "objects": [  
    {  
      "id": "ScheduleId113",  
      "startDateTime": "2013-08-26T00:00:00",  
      "name": "My Copy Schedule",  
      "type": "Schedule",  
      "period": "1 Days"  
    },  
    {  
      "id": "CopyActivityId112",  
      "input": {  
        "ref": " MySqlDataNodeId115"  
      },  
      "schedule": {  
        "ref": "ScheduleId113"  
      },  
      "name": "My Copy",  
      "runsOn": {  
        "ref": "Ec2ResourceId116"  
      },  
      "onSuccess": {  
        "ref": "ActionId1"  
      },  
      "onFail": {  
        "ref": "SnsAlarmId117"  
      },  
      "output": {  
        "ref": "S3DataNodeId114"  
      }  
    }  
  ]  
}
```

```
        },
        "type": "CopyActivity"
    },
    {
        "id": "S3DataNodeId114",
        "schedule": {
            "ref": "ScheduleId113"
        },
        "filePath": "s3://amzn-s3-demo-bucket/rds-output/output.csv",
        "name": "My S3 Data",
        "type": "S3DataNode"
    },
    {
        "id": " MySqlDataNodeId115",
        "username": "my-username",
        "schedule": {
            "ref": "ScheduleId113"
        },
        "name": "My RDS Data",
        "*password": "my-password",
        "table": "table-name",
        "connectionString": "jdbc:mysql://your-sql-instance-name.id.region-name.rds.amazonaws.com:3306/database-name",
        "selectQuery": "select * from #{table}",
        "type": "SqlDataNode"
    },
    {
        "id": "Ec2ResourceId116",
        "schedule": {
            "ref": "ScheduleId113"
        },
        "name": "My EC2 Resource",
        "role": "DataPipelineDefaultRole",
        "type": "Ec2Resource",
        "resourceRole": "DataPipelineDefaultResourceRole"
    },
    {
        "message": "This is a success message.",
        "id": "ActionId1",
        "subject": "RDS to S3 copy succeeded!",
        "name": "My Success Alarm",
        "role": "DataPipelineDefaultRole",
        "topicArn": "arn:aws:sns:us-east-1:123456789012:example-topic",
        "type": "SnsAlarm"
```

```

},
{
  "id": "Default",
  "scheduleType": "timeseries",
  "failureAndRerunMode": "CASCADE",
  "name": "Default",
  "role": "DataPipelineDefaultRole",
  "resourceRole": "DataPipelineDefaultResourceRole"
},
{
  "message": "There was a problem executing #{node.name} at for period
#{node.@scheduledStartTime} to #{node.@scheduledEndTime}",
  "id": "SnsAlarmId117",
  "subject": "RDS to S3 copy failed",
  "name": "My Failure Alarm",
  "role": "DataPipelineDefaultRole",
  "topicArn": "arn:aws:sns:us-east-1:123456789012:example-topic",
  "type": "SnsAlarm"
}
]
}

```

MySQL-Datenknoten

Die MySqlDataNode Eingabe-Pipeline-Komponente definiert einen Speicherort für die Eingabedaten, in diesem Fall eine Amazon RDS-Instance. Die MySqlDataNode Eingabekomponente wird durch die folgenden Felder definiert:

```

{
  "id": "MySqlDataNodeId115",
  "username": "my-username",
  "schedule": {
    "ref": "ScheduleId113"
  },
  "name": "My RDS Data",
  "*password": "my-password",
  "table": "table-name",
  "connectionString": "jdbc:mysql://your-sql-instance-name.id.region-name.rds.amazonaws.com:3306/database-name",
  "selectQuery": "select * from #{table}",
  "type": "SqlDataNode"
},

```

Id

Den benutzerdefinierten Namen, der nur als Referenz dient.

Username

Der Benutzername des Datenbank-Kontos, das über ausreichend Berechtigungen verfügt, um Daten aus der Datenbanktabelle abzurufen. *my-username* Ersetzen Sie durch den Namen Ihres Benutzers.

Plan

Ein Verweis auf die Zeitplankomponente, die wir in den vorhergehenden Zeilen der JSON-Datei erstellt haben.

Name

Den benutzerdefinierten Namen, der nur als Referenz dient.

*Passwort

Das Passwort für das Datenbankkonto mit dem Sternchen-Präfix, das angibt, dass der Kennwortwert verschlüsselt AWS Data Pipeline werden muss. *my-password* Ersetzen Sie es durch das richtige Passwort für Ihren Benutzer. Dem Passwortfeld ist das Sternchen-Sonderzeichen vorangestellt. Weitere Informationen finden Sie unter [Sonderzeichen](#).

Tabelle

Der Name der Datenbanktabelle, die die zu kopierenden Daten enthält. *table-name* Ersetzen Sie es durch den Namen Ihrer Datenbanktabelle.

connectionString

Die JDBC-Verbindungszeichenfolge für das CopyActivity Objekt, das eine Verbindung mit der Datenbank herstellen soll.

selectQuery

Eine gültige SQL-SELECT-Abfrage, die festlegt, welche Daten aus der Datenbanktabelle kopiert werden sollen. Beachten Sie, dass `#{table}` ein Ausdruck ist, der den Tabellennamen wiederverwendet, der in den vorhergehenden Zeilen der JSON-Datei durch die Variable "table" angegeben wird.

Typ

Der SqlDataNode Typ, bei dem es sich in diesem Beispiel um eine Amazon RDS-Instance handelt, die MySQL verwendet.

Note

Der MySqlDataNode Typ ist veraltet. Sie können es zwar weiterhin verwenden MySqlDataNode, wir empfehlen jedoch, es zu verwenden. SqlDataNode

Datenknoten

Als Nächstes definiert die Pipeline-Komponente S3Output einen Speicherort für die Ausgabedatei; in diesem Fall eine CSV-Datei an einem Amazon S3 S3-Bucket-Speicherort. Die DataNode S3-Ausgabekomponente wird durch die folgenden Felder definiert:

```
{  
  "id": "S3DataNodeId114",  
  "schedule": {  
    "ref": "ScheduleId113"  
  },  
  "filePath": "s3://amzn-s3-demo-bucket/rds-output/output.csv",  
  "name": "My S3 Data",  
  "type": "S3DataNode"  
},
```

Id

Die benutzerdefinierte ID, die nur als Referenz dient.

Plan

Ein Verweis auf die Zeitplankomponente, die wir in den vorhergehenden Zeilen der JSON-Datei erstellt haben.

filePath

Den Pfad zu den Daten, die mit dem Datenknoten verknüpft sind, der in diesem Beispiel eine CSV-Ausgabedatei ist.

Name

Den benutzerdefinierten Namen, der nur als Referenz dient.

Typ

Der Pipeline-Objekttyp, der S3 ist, DataNode um dem Ort zu entsprechen, an dem sich die Daten befinden, in einem Amazon S3 S3-Bucket.

Ressource

Dies ist eine Definition der Rechenressource, die die Kopieroperation ausführt. In diesem Beispiel AWS Data Pipeline sollte automatisch eine EC2 Instanz zur Ausführung der Kopieraufgabe erstellt und die Ressource nach Abschluss der Aufgabe beendet werden. Die hier definierten Felder steuern die Erstellung und Funktion der EC2 Instanz, die die Arbeit erledigt. Die EC2 Ressource wird durch die folgenden Felder definiert:

```
{  
  "id": "Ec2ResourceId116",  
  "schedule": {  
    "ref": "ScheduleId113"  
  },  
  "name": "My EC2 Resource",  
  "role": "DataPipelineDefaultRole",  
  "type": "Ec2Resource",  
  "resourceRole": "DataPipelineDefaultResourceRole"  
},
```

Id

Die benutzerdefinierte ID, die nur als Referenz dient.

Plan

Der Zeitplan für die Erstellung dieser Rechenressource.

Name

Den benutzerdefinierten Namen, der nur als Referenz dient.

Rolle

Die IAM-Rolle des Kontos, das auf Ressourcen zugreift, z. B. auf einen Amazon S3 S3-Bucket zum Abrufen von Daten.

Typ

Der Typ der Rechenressource, die Arbeit ausführen soll; in diesem Fall eine Instance. EC2 Es sind auch andere Ressourcentypen verfügbar, z. B. ein EmrCluster Typ.

resourceRole

Die IAM-Rolle des Kontos, das Ressourcen erstellt, z. B. das Erstellen und Konfigurieren einer EC2 Instanz in Ihrem Namen. Rolle und ResourceRole können dieselbe Rolle sein, bieten jedoch getrennt eine größere Granularität in Ihrer Sicherheitskonfiguration.

Aktivität

Der letzte Abschnitt in der JSON-Datei ist die Definition der Aktivität, die die auszuführende Arbeit darstellt. In diesem Fall verwenden wir eine CopyActivity Komponente, um Daten aus einer Datei in einem Amazon-S3-Bucket in eine andere Datei zu kopieren. Die CopyActivity Komponente wird durch die folgenden Felder definiert:

```
{  
    "id": "CopyActivityId112",  
    "input": {  
        "ref": "MySQLDataNode115"  
    },  
    "schedule": {  
        "ref": "ScheduleId113"  
    },  
    "name": "My Copy",  
    "runsOn": {  
        "ref": "Ec2ResourceId116"  
    },  
    "onSuccess": {  
        "ref": "ActionId1"  
    },  
    "onFail": {  
        "ref": "SnsAlarmId117"  
    },  
    "output": {  
        "ref": "S3DataNodeId114"  
    },  
    "type": "CopyActivity"  
},
```

Id

Die benutzerdefinierte ID, die nur als Referenz dient

Eingabe

Den Speicherort der zu kopierenden MySQL-Daten

Plan

Den Zeitplan für die Ausführung dieser Aktivität

Name

Den benutzerdefinierten Namen, der nur als Referenz dient

runsOn

Die Datenverarbeitungsressource, die die Arbeit ausführt, die diese Aktivität definiert. In diesem Beispiel geben wir einen Verweis auf die zuvor definierte EC2 Instanz an. Wenn Sie das `runsOn` Feld verwenden AWS Data Pipeline , wird die EC2 Instanz für Sie erstellt. Das Feld `runsOn` zeigt an, dass die Ressource in der AWS-Infrastruktur vorhanden ist, während der Wert `workerGroup` angibt, dass Sie Ihre eigenen lokalen Ressourcen zur Ausführung der Arbeit verwenden möchten.

onSuccess

Den [SnsAlarm](#), der versendet werden soll, wenn die Aktivität erfolgreich abgeschlossen wurde.

onFail

Den [SnsAlarm](#), der versendet werden soll, wenn die Aktivität fehlschlägt.

Output

Der Amazon-S3-Speicherort der CSV-Ausgabedatei

Typ

Den Typ der Aktivität, die durchgeführt werden soll.

Hochladen und Aktivieren der Pipeline-Definition

Sie müssen Ihre Pipeline-Definition hochladen und Ihre Pipeline aktivieren. Ersetzen Sie die Befehle in den folgenden Beispielbefehlen `pipeline_name` durch ein Label für Ihre Pipeline und `pipeline_file` durch den vollqualifizierten Pfad für die .json Pipeline-Definitionsdatei.

AWS CLI

Verwenden Sie den folgenden Befehl [create-pipeline](#), um Ihre Pipeline-Definition zu erstellen und Ihre Pipeline zu aktivieren. Notieren Sie sich die ID Ihrer Pipeline, da Sie diesen Wert mit den meisten CLI-Befehlen verwenden werden.

```
aws datapipeline create-pipeline --name pipeline_name --unique-id token
{
    "pipelineId": "df-00627471S0VYZEXAMPLE"
}
```

Verwenden Sie den folgenden [put-pipeline-definition](#)Befehl, um Ihre Pipeline-Definition hochzuladen.

```
aws datapipeline put-pipeline-definition --pipeline-id df-00627471S0VYZEXAMPLE --
pipeline-definition file://MyEmrPipelineDefinition.json
```

Wenn Ihre Pipeline erfolgreich validiert wurde, ist das validationErrors Feld leer. Sie sollten alle Warnungen überprüfen.

Verwenden Sie den folgenden Befehl [activate-pipeline](#), um Ihre Pipeline zu aktivieren.

```
aws datapipeline activate-pipeline --pipeline-id df-00627471S0VYZEXAMPLE
```

Mit dem folgenden Befehl [list-pipelines](#) können Sie überprüfen, ob Ihre Pipeline in der Pipeline-Liste aufgeführt ist.

```
aws datapipeline list-pipelines
```

Daten mithilfe von Amazon Redshift nach Amazon Redshift kopieren AWS Data Pipeline

Dieses Tutorial führt Sie durch den Prozess zum Erstellen einer Pipeline, die regelmäßig Daten von Amazon S3 in Amazon Redshift verschiebt, indem Sie entweder die Vorlage Copy to Redshift in Amazon Redshift verschiebt, oder eine Pipeline-Definitionsdatei mit der AWS Data Pipeline CLI. AWS Data Pipeline

Amazon S3 ist ein Webservice, der es Ihnen ermöglicht, Daten in der Cloud zu speichern. Weitere Informationen finden Sie im [Benutzerhandbuch für Amazon Simple Storage Service](#).

Amazon Redshift ist ein Data-Warehouse-Service in der Cloud. Weitere Informationen finden Sie im [Amazon Redshift Redshift-Verwaltungshandbuch](#).

Für dieses Tutorial gibt es mehrere Voraussetzungen. Nachdem Sie die folgenden Schritte ausgeführt haben, können Sie das Tutorial über die Konsole oder die Befehlszeile fortsetzen.

Inhalt

- [Bevor Sie beginnen: Konfigurieren Sie COPY-Optionen und laden Sie Daten](#)
- [Pipeline einrichten, Sicherheitsgruppe erstellen und Amazon Redshift Redshift-Cluster erstellen](#)
- [Daten über die Befehlszeile nach Amazon Redshift kopieren](#)

Bevor Sie beginnen: Konfigurieren Sie COPY-Optionen und laden Sie Daten

Bevor Sie Daten nach Amazon Redshift innerhalb kopieren, stellen Sie sicher AWS Data Pipeline, dass Sie:

- Laden Sie Daten aus Amazon S3
- Richten Sie die COPY Aktivität in Amazon Redshift ein.

Wenn Sie diese Optionen aktiviert und erfolgreich Daten geladen haben, übertragen Sie diese Optionen in AWS Data Pipeline, um den Kopiervorgang darin auszuführen.

COPYOptionen finden Sie unter [COPY](#) im Amazon Redshift Database Developer Guide.

Schritte zum Laden von Daten aus Amazon S3 finden Sie unter [Daten aus Amazon S3 laden](#) im Amazon Redshift Database Developer Guide.

Der folgende SQL-Befehl in Amazon Redshift erstellt beispielsweise eine neue Tabelle mit dem Namen LISTING und kopiert Beispieldaten aus einem öffentlich verfügbaren Bucket in Amazon S3.

Ersetzen Sie den <iam-role-arn> und die Region durch Ihre eigenen Werte.

Einzelheiten zu diesem Beispiel finden Sie unter [Beispieldaten aus Amazon S3 laden im Amazon Redshift Getting Started Guide](#).

```
create table listing(
    listid integer not null distkey,
    sellerid integer not null,
    eventid integer not null,
    dateid smallint not null sortkey,
    numtickets smallint not null,
    priceperticket decimal(8,2),
    totalprice decimal(8,2),
    listtime timestamp);
```

```
copy listing from 's3://awssampledbuswest2/ticket/listings_pipe.txt'
credentials 'aws_iam_role=<iam-role-arn>'
delimiter '|' region 'us-west-2';
```

Pipeline einrichten, Sicherheitsgruppe erstellen und Amazon Redshift Redshift-Cluster erstellen

So richten Sie Ihr System für das Tutorial ein

1. Führen Sie die Aufgaben unter au [Einrichtung für AWS Data Pipeline](#).
2. Erstellen einer Sicherheitsgruppe.
 - a. Öffnen Sie die EC2 Amazon-Konsole.
 - b. Klicken Sie im Navigationsbereich auf Security Groups.
 - c. Klicken Sie auf Create Security Group.
 - d. Geben Sie einen Namen und eine Beschreibung für die Sicherheitsgruppe an.
 - e. [EC2-Classic] Wählen Sie No VPC für VPC.
 - f. [EC2-VPC] Wählen Sie die ID Ihrer VPC für VPC aus.
 - g. Klicken Sie auf Create.
3. [EC2-Classic] Erstellen Sie eine Amazon Redshift Redshift-Cluster-Sicherheitsgruppe und geben Sie die EC2 Amazon-Sicherheitsgruppe an.
 - a. Öffnen Sie die Amazon-Redshift-Konsole.
 - b. Klicken Sie im Navigationsbereich auf Security Groups.
 - c. Klicken Sie auf Create Cluster Security Group.
 - d. Geben Sie im Dialogfeld Create Cluster Security Group einen Namen und eine Beschreibung für die Cluster-Sicherheitsgruppe an.
 - e. Klicken Sie auf den Namen der neuen Cluster-Sicherheitsgruppe.
 - f. Klicken Sie auf Add Connection Type.
 - g. Wählen Sie im Dialogfeld Verbindungstyp hinzufügen unter Verbindungstyp die Option EC2 Sicherheitsgruppe aus, wählen Sie unter Sicherheitsgruppenname die Sicherheitsgruppe ausEC2 , die Sie erstellt haben, und klicken Sie dann auf Autorisieren.
4. [EC2-VPC] Erstellen Sie eine Amazon Redshift Redshift-Cluster-Sicherheitsgruppe und geben Sie die VPC-Sicherheitsgruppe an.

- a. Öffnen Sie die EC2 Amazon-Konsole.
 - b. Klicken Sie im Navigationsbereich auf Security Groups.
 - c. Klicken Sie auf Create Security Group.
 - d. Geben Sie im Dialogfeld Create Security Group einen Namen und eine Beschreibung für die Sicherheitsgruppe an und wählen Sie für VPC die ID Ihres VPC aus.
 - e. Klicken Sie auf Add Rule. Geben Sie den Typ, das Protokoll und den Portbereich an und geben Sie die ID der Sicherheitsgruppe in Source ein. Wählen Sie die Sicherheitsgruppe aus, die Sie im zweiten Schritt erstellt haben.
 - f. Klicken Sie auf Create.
5. Nachfolgend finden Sie eine kurze Zusammenfassung der Schritte.

Wenn Sie bereits über einen Amazon-Redshift-Cluster verfügen, notieren Sie sich die Cluster-ID.

Um einen neuen Cluster zu erstellen und Beispieldaten zu laden, folgen Sie den Schritten unter [Erste Schritte mit Amazon Redshift](#). Weitere Informationen zum Erstellen von Clustern finden Sie unter [Creating a Cluster](#) im Amazon Redshift Management Guide.

- a. Öffnen Sie die Amazon-Redshift-Konsole.
- b. Klicken Sie auf Launch Cluster.
- c. Geben Sie die erforderlichen Details für Ihren Cluster an und klicken Sie dann auf Continue.
- d. Geben Sie die Knotenkonfiguration an und klicken Sie dann auf Continue.
- e. Wählen Sie auf der Seite für zusätzliche Konfigurationsdaten die von Ihnen erstellte Cluster-Sicherheitsgruppe aus und klicken Sie dann auf Continue.
- f. Überprüfen Sie die Spezifikationen für Ihren Cluster und klicken Sie dann auf Launch Cluster.

Daten über die Befehlszeile nach Amazon Redshift kopieren

Dieses Tutorial zeigt, wie Sie Daten aus Amazon S3 in Amazon Redshift kopieren. Sie erstellen eine neue Tabelle in Amazon Redshift und verwenden sie dann, AWS Data Pipeline um Daten aus einem öffentlichen Amazon S3 S3-Bucket, der Beispieleingabedaten im CSV-Format enthält, in diese Tabelle zu übertragen. Die Protokolle werden in einem Amazon-S3-Bucket gespeichert, den Sie besitzen.

Amazon S3 ist ein Webservice, der es Ihnen ermöglicht, Daten in der Cloud zu speichern. Weitere Informationen finden Sie im [Benutzerhandbuch für Amazon Simple Storage Service](#). Amazon Redshift ist ein Data-Warehouse-Service in der Cloud. Weitere Informationen finden Sie im [Amazon Redshift Redshift-Verwaltungshandbuch](#).

Voraussetzungen

Bevor Sie beginnen, müssen Sie die folgenden Schritte ausführen:

1. Installieren und konfigurieren Sie eine Befehlszeilenschnittstelle (CLI). Weitere Informationen finden Sie unter [Zugreifen AWS Data Pipeline](#).
2. Stellen Sie sicher, dass die IAM-Rollen benannt DataPipelineDefaultResourceRolesind DataPipelineDefaultRoleund existieren. Die AWS Data Pipeline Konsole erstellt diese Rollen automatisch für Sie. Wenn Sie die AWS Data Pipeline Konsole nicht mindestens einmal verwendet haben, müssen Sie diese Rollen manuell erstellen. Weitere Informationen finden Sie unter [IAM-Rollen für AWS Data Pipeline](#).
3. Richten Sie den COPY Befehl in Amazon Redshift ein, da dieselben Optionen funktionieren müssen, wenn Sie das Kopieren innerhalb von AWS Data Pipeline Amazon Redshift durchführen. Weitere Informationen finden Sie unter [Bevor Sie beginnen: Konfigurieren Sie COPY-Optionen und laden Sie Daten](#).
4. Einrichten einer Amazon-Redshift-Datenbank. Weitere Informationen finden Sie unter [Pipeline einrichten, Sicherheitsgruppe erstellen und Amazon Redshift Redshift-Cluster erstellen](#).

Aufgaben

- [Definieren Sie eine Pipeline im JSON-Format](#)
- [Hochladen und Aktivieren der Pipeline-Definition](#)

Definieren Sie eine Pipeline im JSON-Format

Dieses Beispielszenario zeigt, wie Daten aus einem Amazon-S3-Bucket in Amazon Redshift kopiert werden.

Dies ist die vollständige Pipeline-Definition-JSON-Datei, gefolgt von einer Erläuterung für jeden ihrer Abschnitte. Wir empfehlen, dass Sie einen Texteditor verwenden, mit dem Sie die Syntax von JSON-formatierten Dateien überprüfen können, und die Datei mit der Dateierweiterung .json benennen.

{

```
"objects": [
  {
    "id": "CSVID1",
    "name": "DefaultCSV1",
    "type": "CSV"
  },
  {
    "id": "RedshiftDatabaseId1",
    "databaseName": "dbname",
    "username": "user",
    "name": "DefaultRedshiftDatabase1",
    "*password": "password",
    "type": "RedshiftDatabase",
    "clusterId": "redshiftclusterId"
  },
  {
    "id": "Default",
    "scheduleType": "timeseries",
    "failureAndRerunMode": "CASCADE",
    "name": "Default",
    "role": "DataPipelineDefaultRole",
    "resourceRole": "DataPipelineDefaultResourceRole"
  },
  {
    "id": "RedshiftDataNodeId1",
    "schedule": {
      "ref": "ScheduleId1"
    },
    "tableName": "orders",
    "name": "DefaultRedshiftDataNode1",
    "createTableSql": "create table StructuredLogs (requestBeginTime CHAR(30)
PRIMARY KEY DISTKEY SORTKEY, requestEndTime CHAR(30), hostname CHAR(100), requestDate
varchar(20));",
    "type": "RedshiftDataNode",
    "database": {
      "ref": "RedshiftDatabaseId1"
    }
  },
  {
    "id": "Ec2ResourceId1",
    "schedule": {
      "ref": "ScheduleId1"
    },
    "securityGroups": "MySecurityGroup",
```

```
"name": "DefaultEc2Resource1",
"role": "DataPipelineDefaultRole",
"logUri": "s3://myLogs",
"resourceRole": "DataPipelineDefaultResourceRole",
"type": "Ec2Resource"
},
{
  "id": "ScheduleId1",
  "startDateTime": "yyyy-mm-ddT00:00:00",
  "name": "DefaultSchedule1",
  "type": "Schedule",
  "period": "period",
  "endDateTime": "yyyy-mm-ddT00:00:00"
},
{
  "id": "S3DataNodeId1",
  "schedule": {
    "ref": "ScheduleId1"
  },
  "filePath": "s3://datapipeline-us-east-1/samples/hive-ads-samples.csv",
  "name": "DefaultS3DataNode1",
  "dataFormat": {
    "ref": "CSVId1"
  },
  "type": "S3DataNode"
},
{
  "id": "RedshiftCopyActivityId1",
  "input": {
    "ref": "S3DataNodeId1"
  },
  "schedule": {
    "ref": "ScheduleId1"
  },
  "insertMode": "KEEP_EXISTING",
  "name": "DefaultRedshiftCopyActivity1",
  "runsOn": {
    "ref": "Ec2ResourceId1"
  },
  "type": "RedshiftCopyActivity",
  "output": {
    "ref": "RedshiftDataNodeId1"
  }
}
```

```
]  
}
```

Weitere Informationen zu diesen Objekten finden Sie in der folgenden Dokumentation.

Objekte

- [Datenknoten](#)
- [Ressource](#)
- [Aktivität](#)

Datenknoten

Bei diesem Beispiel werden ein Eingabedatenknoten, ein Ausgabedatenknoten und eine Datenbank verwendet.

Eingabedatenknoten

Die S3DataNode Eingabe-Pipeline-Komponente definiert den Speicherort der Eingabedaten in Amazon S3 und das Datenformat der Eingabedaten. Weitere Informationen finden Sie unter [S3 DataNode](#).

Diese Eingabekomponente wird durch folgende Felder definiert:

```
{  
  "id": "S3DataNodeId1",  
  "schedule": {  
    "ref": "ScheduleId1"  
  },  
  "filePath": "s3://datapipeline-us-east-1/samples/hive-ads-samples.csv",  
  "name": "DefaultS3DataNode1",  
  "dataFormat": {  
    "ref": "CSVID1"  
  },  
  "type": "S3DataNode"  
},
```

id

Die benutzerdefinierte ID, die nur als Referenz dient.

schedule

Einen Verweis auf die Zeitplankomponente.

filePath

Den Pfad zu den Daten, die mit dem Datenknoten verknüpft sind, der in diesem Beispiel eine CSV-Eingabedatei ist.

name

Den benutzerdefinierten Namen, der nur als Referenz dient.

dataFormat

Einen Verweis auf das Format der Daten für die Aktivitätsverarbeitung.

Ausgabedatenknoten

Die RedshiftDataNode Ausgabe-Pipeline-Komponente definiert einen Speicherort für die Ausgabedaten, in diesem Fall eine Tabelle in einer Amazon Redshift Redshift-Datenbank. Weitere Informationen finden Sie unter [RedshiftDataNode](#). Diese Ausgabekomponente wird durch folgende Felder definiert:

```
{  
  "id": "RedshiftDataNodeId1",  
  "schedule": {  
    "ref": "ScheduleId1"  
  },  
  "tableName": "orders",  
  "name": "DefaultRedshiftDataNode1",  
  "createTableSql": "create table StructuredLogs (requestBeginTime CHAR(30) PRIMARY  
KEY DISTKEY SORTKEY, requestEndTime CHAR(30), hostname CHAR(100), requestDate  
varchar(20));",  
  "type": "RedshiftDataNode",  
  "database": {  
    "ref": "RedshiftDatabaseId1"  
  }  
},
```

id

Die benutzerdefinierte ID, die nur als Referenz dient.

schedule

Einen Verweis auf die Zeitplankomponente.

tableName

Der Name der Amazon Redshift-Tabelle.

name

Den benutzerdefinierten Namen, der nur als Referenz dient.

createTableSql

Einen SQL-Ausdruck, der die Tabelle in der Datenbank erstellt.

database

Ein Verweis auf die Amazon-Redshift-Datenbank.

Datenbank

Die RedshiftDatabase-Komponente wird durch folgende Felder definiert: Weitere Informationen finden Sie unter [RedshiftDatabase](#).

```
{
  "id": "RedshiftDatabaseId1",
  "databaseName": "dbname",
  "username": "user",
  "name": "DefaultRedshiftDatabase1",
  "*password": "password",
  "type": "RedshiftDatabase",
  "clusterId": "redshiftclusterId"
},
```

id

Die benutzerdefinierte ID, die nur als Referenz dient.

databaseName

Den Namen der logischen Datenbank.

username

Den Benutzernamen für die Verbindung zur Datenbank.

name

Den benutzerdefinierten Namen, der nur als Referenz dient.

password

Das Passwort für die Verbindung zur Datenbank.

clusterId

Die ID des Redshift-Clusters.

Ressource

Dies ist eine Definition der Rechenressource, die die Kopieroperation ausführt. In diesem Beispiel AWS Data Pipeline sollte automatisch eine EC2 Instance zur Ausführung der Kopieraufgabe erstellt und die Instance nach Abschluss der Aufgabe beendet werden. Die hier definierten Felder steuern die Erstellung und Funktion der Instance, die die Arbeit erledigt. Weitere Informationen finden Sie unter [Ec2Resource](#).

Die Ec2Resource wird durch folgende Felder definiert:

```
{  
  "id": "Ec2ResourceId1",  
  "schedule": {  
    "ref": "ScheduleId1"  
  },  
  "securityGroups": "MySecurityGroup",  
  "name": "DefaultEc2Resource1",  
  "role": "DataPipelineDefaultRole",  
  "logUri": "s3://myLogs",  
  "resourceRole": "DataPipelineDefaultResourceRole",  
  "type": "Ec2Resource"  
},
```

id

Die benutzerdefinierte ID, die nur als Referenz dient.

schedule

Der Zeitplan für die Erstellung dieser Rechenressource.

securityGroups

Die Sicherheitsgruppe, die für die Instances im Ressourcenpool verwendet werden soll.

name

Den benutzerdefinierten Namen, der nur als Referenz dient.

role

Die IAM-Rolle des Kontos, das auf Ressourcen zugreift, z. B. auf einen Amazon S3 S3-Bucket zum Abrufen von Daten.

logUri

Der Amazon S3 S3-Zielpfad zum Sichern von Task Runner-Protokollen vonEc2Resource.

resourceRole

Die IAM-Rolle des Kontos, das Ressourcen erstellt, z. B. das Erstellen und Konfigurieren einer EC2 Instance in Ihrem Namen. Rolle und ResourceRole können dieselbe Rolle sein, bieten jedoch getrennt eine größere Granularität in Ihrer Sicherheitskonfiguration.

Aktivität

Der letzte Abschnitt in der JSON-Datei ist die Definition der Aktivität, die die auszuführende Arbeit darstellt. In diesem Fall verwenden wir eine RedshiftCopyActivity Komponente, um Daten aus Amazon S3 in Amazon Redshift zu kopieren. Weitere Informationen finden Sie unter [RedshiftCopyActivity](#).

Die RedshiftCopyActivity-Komponente ist durch folgende Felder definiert:

```
{
  "id": "RedshiftCopyActivityId1",
  "input": {
    "ref": "S3DataNodeId1"
  },
  "schedule": {
    "ref": "ScheduleId1"
  },
  "insertMode": "KEEP_EXISTING",
  "name": "DefaultRedshiftCopyActivity1",
  "runsOn": {
    "ref": "Ec2ResourceId1"
```

```
},
  "type": "RedshiftCopyActivity",
  "output": {
    "ref": "RedshiftDataNodeId1"
  }
},
```

id

Die benutzerdefinierte ID, die nur als Referenz dient.

input

Ein Verweis auf die Amazon-S3-Quelldatei.

schedule

Der Zeitplan für die Ausführung dieser Aktivität.

insertMode

Der Einfügetyp (KEEP_EXISTING, OVERWRITE_EXISTING oder TRUNCATE).

name

Den benutzerdefinierten Namen, der nur als Referenz dient.

runsOn

Die Datenverarbeitungsressource, die die Arbeit ausführt, die diese Aktivität definiert.

output

Ein Verweis auf die Amazon Redshift Redshift-Zieltablette.

Hochladen und Aktivieren der Pipeline-Definition

Sie müssen Ihre Pipeline-Definition hochladen und Ihre Pipeline aktivieren. Ersetzen Sie die Befehle in den folgenden Beispielbefehlen *pipeline_name* durch ein Label für Ihre Pipeline und *pipeline_file* durch den vollqualifizierten Pfad für die .json Pipeline-Definitionsdatei.

AWS CLI

Verwenden Sie den folgenden Befehl [create-pipeline](#), um Ihre Pipeline-Definition zu erstellen und Ihre Pipeline zu aktivieren. Notieren Sie sich die ID Ihrer Pipeline, da Sie diesen Wert mit den meisten CLI-Befehlen verwenden werden.

```
aws datapipeline create-pipeline --name pipeline_name --unique-id token
{
    "pipelineId": "df-00627471S0VYZEXAMPLE"
}
```

Verwenden Sie den folgenden [put-pipeline-definition](#)Befehl, um Ihre Pipeline-Definition hochzuladen.

```
aws datapipeline put-pipeline-definition --pipeline-id df-00627471S0VYZEXAMPLE --
pipeline-definition file://MyEmrPipelineDefinition.json
```

Wenn Ihre Pipeline erfolgreich validiert wurde, ist das validationErrors Feld leer. Sie sollten alle Warnungen überprüfen.

Verwenden Sie den folgenden Befehl [activate-pipeline](#), um Ihre Pipeline zu aktivieren.

```
aws datapipeline activate-pipeline --pipeline-id df-00627471S0VYZEXAMPLE
```

Mit dem folgenden Befehl [list-pipelines](#) können Sie überprüfen, ob Ihre Pipeline in der Pipeline-Liste aufgeführt ist.

```
aws datapipeline list-pipelines
```

Pipeline-Ausdrücke und -Funktionen

Dieser Abschnitt enthält Informationen zur Syntax der in Pipelines verwendeten Ausdrücke und Funktionen sowie zu den zugehörigen Datentypen.

Einfache Datentypen

Die folgenden Datentypen sind als Feldwerte zulässig.

Typen

- [DateTime](#)
- [Numerischer Wert](#)
- [Objektverweise](#)
- [Intervall](#)
- [String](#)

DateTime

AWS Data Pipeline unterstützt das Datum und die im Format „J-MM-TTSS:MM:SS“ nur in UTC/GMT. Im folgenden Beispiel wird dem Feld `startDateTime` eines Schedule-Objekts der Wert `1/15/2012, 11:59 p.m.` in der UTC/GMT-Zeitzone zugewiesen.

```
"startDateTime" : "2012-01-15T23:59:00"
```

Numerischer Wert

AWS Data Pipeline unterstützt sowohl Ganzzahlen als auch Gleitkommawerte.

Objektverweise

Dies ist ein Objekt in der Pipeline-Definition. Dabei kann es sich um das aktuelle Objekt, um den Namen eines woanders in der Pipeline definierten Objekts oder um ein Objekt mit einem Feld handeln, in dem mit dem Schlüsselwort `node` auf das aktuelle Objekt verwiesen wird. Mehr über `node` erfahren Sie unter [Verweisen auf Felder und Objekte](#). Weitere Informationen zu den Pipeline-Objekttypen finden Sie unter [Pipeline-Objektreferenz](#).

Intervall

Gibt an, wie oft ein geplantes Ereignis ausgeführt werden soll. Die Angabe erfolgt im Format "N [years|months|weeks|days|hours|minutes]", wobei N eine positive Ganzzahl ist.

Der Mindestzeitraum beträgt 15 Minuten und der maximale Zeitraum beträgt 3 Jahre.

Im folgenden Beispiel wird das Feld `period` des Objekts `Schedule` auf 3 Stunden eingestellt. Dadurch wird ein geplantes Ereignis alle drei Stunden ausgeführt.

```
"period" : "3 hours"
```

String

Standard-Zeichenfolgenwerte. Zeichenfolgen müssen in Anführungszeichen ("") eingeschlossen werden. Der umgekehrte Schrägstrich (\) kann als Escape-Zeichen in einer Zeichenfolge verwendet werden. Mehrzeilige Zeichenfolgen werden nicht unterstützt.

Die folgenden Beispiele sind gültige Zeichenfolgenwerte für das Feld `id`.

```
"id" : "My Data Object"  
"id" : "My \"Data\" Object"
```

Zeichenfolgen können auch Ausdrücke enthalten, die zu Zeichenfolgenwerten ausgewertet werden. Diese Ausdrücke müssen in der Zeichenfolge zwischen den Trennzeichen "#" und "}" stehen. Im folgenden Beispiel wird mit einem Ausdruck der Name des aktuellen Objekts in einen Pfad eingefügt.

```
"filePath" : "s3://amzn-s3-demo-bucket/#${name}.csv"
```

Weitere Informationen zur Arbeit mit Ausdrücken finden Sie unter [Verweisen auf Felder und Objekte](#) und [Ausdrucksauswertung](#).

Ausdrücke

Ausdrücke ermöglichen die Nutzung eines Werts in mehreren zusammengehörigen Objekten. Ausdrücke werden zur Laufzeit vom AWS Data Pipeline Webdienst verarbeitet, wodurch sichergestellt wird, dass alle Ausdrücke durch den Wert des Ausdrucks ersetzt werden.

Ausdrücke müssen in die "#{" und "}" Trennzeichen eingeschlossen werden. Ein Ausdruck kann in jedem Pipeline-Definitionsobjekt verwendet werden, in dem Zeichenfolgen zulässig sind. Wenn ein Slot ein Verweis ist bzw. den Typ ID, NAME, TYPE oder SPHERE hat, wird sein Wert nicht ausgewertet und er wird unverändert übernommen.

Der folgende Ausdruck ruft eine der AWS Data Pipeline Funktionen auf. Weitere Informationen finden Sie unter [Ausdrucksauswertung](#).

```
#{{format(myDateTime, 'YYYY-MM-dd hh:mm:ss')}}
```

Verweisen auf Felder und Objekte

Ausdrücke können Felder des aktuellen Objekts, in dem sie sich befinden, oder Felder eines anderen Objekts verwenden, das durch einen Verweis verknüpft ist.

Ein Slot-Format besteht aus einem Zeitpunkt der Erstellung, gefolgt von dem Zeitpunkt der Objekterstellung, z. B. @S3BackupLocation_2018-01-31T11:05:33.

Sie können auch auf die genaue Steckplatz-ID verweisen, die in der Pipeline-Definition angegeben ist, z. B. die Steckplatz-ID des Amazon S3 S3-Backup-Speicherorts. Um auf die Slot-ID zu verweisen, verwenden Sie #{parent.@id}.

Im folgenden Beispiel verweist das Feld filePath auf das Feld id desselben Objekts, um einen Dateinamen zu bilden. Der Wert von filePath ergibt sich als „s3://amzn-s3-demo-bucket/ExampleDataNode.csv“.

```
{
  "id" : "ExampleDataNode",
  "type" : "S3DataNode",
  "schedule" : {"ref" : "ExampleSchedule"},
  "filePath" : "s3://amzn-s3-demo-bucket/#{parent.@id}.csv",
  "precondition" : {"ref" : "ExampleCondition"},
  "onFail" : {"ref" : "FailureNotify"}
}
```

Um auf das Feld eines anderen, durch einen Verweis verknüpften Objekts zuzugreifen, muss das Schlüsselwort node angegeben werden. Dieses Schlüsselwort ist nur für Alarm- und Vorbedingungsobjekte verfügbar.

Zurück zum obigen Beispiel: Ein Ausdruck in einem SnsAlarm-Objekt kann auf den Datums- und Uhrzeitbereich in einem Schedule-Objekt zugreifen, da das S3DataNode-Objekt auf beide verweist.

Insbesondere kann im FailureNotify-Objekt das Feld message die Laufzeitfelder @scheduledStartTime und @scheduledEndTime von ExampleSchedule verwenden, da im ExampleDataNode-Objekt das Feld onFail auf FailureNotify und das Feld schedule auf ExampleSchedule verweist.

```
{  
    "id" : "FailureNotify",  
    "type" : "SnsAlarm",  
    "subject" : "Failed to run pipeline component",  
    "message": "Error for interval  
#{node.@scheduledStartTime}..#{node.@scheduledEndTime}.",  
    "topicArn": "arn:aws:sns:us-east-1:28619EXAMPLE:ExampleTopic"  
},
```

Note

Sie können Pipelines erstellen, die Abhängigkeiten enthalten, wie z. B. Pipeline-Aufgaben, die von der Verarbeitung anderer Systeme oder Aufgaben abhängen. Wenn Ihre Pipeline bestimmte Ressourcen erfordert, fügen Sie diese Abhängigkeiten mithilfe von Vorbedingungen hinzu, die Sie dann den gewünschten Datenknoten und Aufgaben zuordnen. Die Pipeline lässt sich dann einfacher debuggen und ist weniger fehleranfällig. Verwenden Sie die Abhängigkeiten möglichst in einer einzigen Pipeline, da die Pipeline-übergreifende Fehlersuche schwierig ist.

Verschachtelte Ausdrücke

AWS Data Pipeline ermöglicht es Ihnen, Werte zu verschachteln, um komplexere Ausdrücke zu erstellen. Um beispielsweise eine Zeitberechnung durchzuführen (30 Minuten von scheduledStartTime subtrahieren) und das Ergebnis für die Verwendung in einer Pipeline-Definition zu formatieren, könnten Sie folgenden Ausdruck in einer Aktivität benutzen

```
#{{format(minusMinutes(@scheduledStartTime,30), 'YYYY-MM-dd hh:mm:ss')}}
```

und verwendet das node Präfix, wenn der Ausdruck Teil einer SnsAlarm Oder-Vorbedingung ist:

```
#{{format(minusMinutes(node.@scheduledStartTime,30), 'YYYY-MM-dd hh:mm:ss')}}
```

Listen

In Ausdrücken können auch Listen verwendet werden. Gehen wir von dieser Listendefinition aus:
"myList": ["one", "two"]. Wenn Sie diese Liste nun im Ausdruck #{"this is ' + myList} verwenden, wird sie zu ["this is one", "this is two"] ausgewertet. Wenn zwei Listen vorhanden sind, werden diese von Data Pipeline bei der Auswertung "abgeflacht" (flattened). Ist beispielsweise myList1 als [1, 2] definiert und myList2 als [3, 4], dann wird der Ausdruck #[{myList1}, #{myList2}] zu [1, 2, 3, 4] ausgewertet.

Knotenausdruck

AWS Data Pipeline verwendet den #{node.*} Ausdruck entweder in SnsAlarm oder PreCondition für einen Rückverweis auf das übergeordnete Objekt einer Pipeline-Komponente. Da der Verweis auf SnsAlarm und PreCondition in einer Aktivität oder Ressource ohne Rückverweis aus diesen Objekten erfolgt, bietet node die Möglichkeit, das verweisende Objekt zu referenzieren. In der folgenden Pipeline-Definition wird beispielsweise gezeigt, wie in einer Fehlerbenachrichtigung mit node auf deren übergeordnetes Objekt (ShellCommandActivity) verwiesen und wie dessen geplante Start- und Endzeit in die SnsAlarm-Nachricht eingefügt werden. Für den scheduledStartTime Verweis auf ShellCommandActivity ist kein node Präfix erforderlich, da er scheduledStartTime sich auf sich selbst bezieht.

 Note

Die mit dem Zeichen "@" beginnenden Felder sind Laufzeitfelder.

```
{  
    "id" : "ShellOut",  
    "type" : "ShellCommandActivity",  
    "input" : {"ref" : "HourlyData"},  
    "command" : "/home/userName/xxx.sh #@scheduledStartTime} #@scheduledEndTime}",  
    "schedule" : {"ref" : "HourlyPeriod"},  
    "stderr" : "/tmp/stderr:#{@scheduledStartTime}",  
    "stdout" : "/tmp/stdout:#{@scheduledStartTime}",  
    "onFail" : {"ref" : "FailureNotify"},  
},  
{  
    "id" : "FailureNotify",  
    "type" : "SnsAlarm",  
    "subject" : "Failed to run pipeline component",
```

```
"message": "Error for interval  
#{node.@scheduledStartTime}..#{node.@scheduledEndTime}.",  
"topicArn":"arn:aws:sns:us-east-1:28619EXAMPLE:ExampleTopic"  
},
```

AWS Data Pipeline unterstützt transitive Verweise für benutzerdefinierte Felder, aber keine Laufzeitfelder. Ein transitiver Verweis ist ein Verweis zwischen zwei Pipeline-Komponenten, die von einer anderen als "Vermittler" dienenden Pipeline-Komponente abhängen. Das folgende Beispiel zeigt einen Verweis auf ein transitives benutzerdefiniertes Feld und einen Verweis auf ein nichttransitives Laufzeitfeld. Beide Verweise sind zulässig. Weitere Informationen finden Sie unter [Benutzerdefinierte Felder](#).

```
{  
  "name": "DefaultActivity1",  
  "type": "CopyActivity",  
  "schedule": {"ref": "Once"},  
  "input": {"ref": "s3nodeOne"},  
  "onSuccess": {"ref": "action"},  
  "workerGroup": "test",  
  "output": {"ref": "s3nodeTwo"}  
},  
{  
  "name": "action",  
  "type": "SnsAlarm",  
  "message": "S3 bucket '#{node.output.directoryPath}' succeeded at  
#{node.@actualEndTime}.",  
  "subject": "Testing",  
  "topicArn": "arn:aws:sns:us-east-1:28619EXAMPLE:ExampleTopic",  
  "role": "DataPipelineDefaultRole"  
}
```

Ausdrucksauswertung

AWS Data Pipeline stellt eine Reihe von Funktionen bereit, mit denen Sie den Wert eines Felds berechnen können. Im folgenden Beispiel wird die Funktion `makeDate` verwendet, um dem Feld `startDate` eines `Schedule`-Objekts die GMT/UTC-Zeit "2011-05-24T0:00:00" zuzuweisen.

```
"startDate" : "makeDate(2011,5,24)"
```

Mathematische Funktionen

Die folgenden Funktionen sind für die Arbeit mit numerischen Werten verfügbar.

Funktion	Beschreibung
+	Addition. Beispiel: #{1 + 2} Ergebnis: 3
-	Subtraktion. Beispiel: #{1 - 2} Ergebnis: -1
*	Multiplikation. Beispiel: #{1 * 2} Ergebnis: 2
/	Division. Wenn Sie zwei Ganzzahlen dividieren, werden die Nachkommastellen abgeschnitten. Beispiel: #{1 / 2}, Ergebnis: 0 Beispiel: #{1.0 / 2}, Ergebnis: .5
^	Exponent. Beispiel: #{2 ^ 2} Ergebnis: 4.0

Funktionen für Zeichenfolgen

Die folgenden Funktionen sind für die Arbeit mit Zeichenfolgen verfügbar.

Funktion	Beschreibung
+	<p>Verkettung. Werte mit einem anderen Typ werden zuerst in Zeichenfolgen konvertiert.</p> <p>Beispiel: #{"hel" + "lo"}</p> <p>Ergebnis: "hello"</p>

Datums- und Zeitfunktionen

Die folgenden Funktionen sind für die Arbeit mit DateTime Werten verfügbar. In den Beispielen hat myDateTime den Wert May 24, 2011 @ 5:10 pm GMT.

Note

Das Datums-/Uhrzeitformat für AWS Data Pipeline ist Joda Time, ein Ersatz für die Java-Datums- und Uhrzeitklassen. Weitere Informationen finden Sie unter [Joda Time - Class. DateFormat](#)

Funktion	Beschreibung
<code>int day(DateTime myDateTime)</code>	<p>Ruft den Tag des DateTime Werts als Ganzzahl ab.</p> <p>Beispiel: #{day(myDateTime)}</p> <p>Ergebnis: 24</p>
<code>int dayOfYear(DateTime myDateTime)</code>	<p>Ruft den Tag des Jahres des DateTime Werts als Ganzzahl ab.</p> <p>Beispiel: #{dayOfYear(myDateTime)}</p>

Funktion	Beschreibung
	Ergebnis: 144
<code>DateTime firstOfMonth(DateTime myDateTime)</code>	<p>Erstellt ein DateTime Objekt für den Anfang des Monats im angegebenen Zeitraum DateTime.</p> <p>Beispiel: <code>#{firstOfMonth(myDateTime)}</code></p> <p>Ergebnis: "2011-05-01T17:10:00z"</p>
<code>String format(DateTime myDateTime, String format)</code>	<p>Erstellt ein String-Objekt, das das Ergebnis der Konvertierung des angegebenen Objekts DateTime unter Verwendung der angegebenen Formatzeichenfolge ist.</p> <p>Beispiel: <code>#{format(myDateTime, 'YYYY-MM-dd HH:mm:ss z')}</code></p> <p>Ergebnis: "2011-05-24T17:10:00 UTC"</p>
<code>int hour(DateTime myDateTime)</code>	<p>Ruft die Stunde des DateTime Werts als Ganzzahl ab.</p> <p>Beispiel: <code>#{hour(myDateTime)}</code></p> <p>Ergebnis: 17</p>

Funktion	Beschreibung
<code>DateTime makeDate(int year,int month,int day)</code>	<p>Erstellt um Mitternacht ein DateTime Objekt in UTC mit dem angegebenen Jahr, Monat und Tag.</p> <p>Beispiel: <code>#{{makeDate(2011,5,24)}}</code></p> <p>Ergebnis: "2011-05-24T0:00:00z"</p>
<code>DateTime makeDateTime(int year,int month,int day,int hour,int minute)</code>	<p>Erstellt ein DateTime Objekt in UTC mit dem angegebenen Jahr, Monat, Tag, Stunde und Minute.</p> <p>Beispiel: <code>#{{makeDateTime(2011,5,24,14,21)}}</code></p> <p>Ergebnis: "2011-05-24T14:21:00z"</p>
<code>DateTime midnight(DateTime myDateTime)</code>	<p>Erstellt ein DateTime Objekt für die aktuelle Mitternacht relativ zur angegebenen Mitternacht DateTime. Hat beispielsweise MyDateTime den Wert 2011-05-25T17:10:00z , lautet das Ergebnis wie folgt.</p> <p>Beispiel: <code>#{{midnight(myDateTime)}}</code></p> <p>Ergebnis: "2011-05-25T0:00:00z"</p>

Funktion	Beschreibung
<code>DateTime minusDays(DateTime myDateTime, int daysToSub)</code>	<p>Erstellt ein DateTime Objekt, das das Ergebnis der Subtraktion der angegebenen Anzahl von Tagen von der angegebenen Anzahl ist. DateTime</p> <p>Beispiel: <code>#{{minusDays(myDateTime,1)}}</code></p> <p>Ergebnis: "2011-05-23T17:10:00z"</p>
<code>DateTime minusHours(DateTime myDateTime, int hoursToSub)</code>	<p>Erstellt ein DateTime Objekt, das das Ergebnis der Subtraktion der angegebenen Anzahl von Stunden von der angegebenen Anzahl ist. DateTime</p> <p>Beispiel: <code>#{{minusHours(myDateTime,1)}}</code></p> <p>Ergebnis: "2011-05-24T16:10:00z"</p>

Funktion	Beschreibung
<code>DateTime minusMinutes(DateTime myDateTime, int minutesToSub)</code>	<p>Erstellt ein DateTime Objekt, das das Ergebnis der Subtraktion der angegebenen Anzahl von Minuten von der angegebenen Anzahl ist. DateTime</p> <p>Beispiel: <code>#{{minusMinutes(myDateTime, 1)}}</code></p> <p>Ergebnis: "2011-05-24T17:09:00z"</p>
<code>DateTime minusMonths(DateTime myDateTime, int monthsToSub)</code>	<p>Erstellt ein DateTime Objekt, das das Ergebnis der Subtraktion der angegebenen Anzahl von Monaten von der angegebenen Anzahl ist. DateTime</p> <p>Beispiel: <code>#{{minusMonths(myDateTime, 1)}}</code></p> <p>Ergebnis: "2011-04-24T17:10:00z"</p>

Funktion	Beschreibung
<code>DateTime minusWeeks(DateTime myDateTime, int weeksToSub)</code>	<p>Erstellt ein DateTime Objekt, das das Ergebnis der Subtraktion der angegebenen Anzahl von Wochen von der angegebenen Anzahl ist. DateTime</p> <p>Beispiel: <code>#{minusWeeks(myDateTime,1)}</code></p> <p>Ergebnis: "2011-05-17T17:10:00z"</p>
<code>DateTime minusYears(DateTime myDateTime, int yearsToSub)</code>	<p>Erstellt ein DateTime Objekt, das sich aus der Subtraktion der angegebenen Anzahl von Jahren von der angegebenen Anzahl von Jahren ergibt. DateTime</p> <p>Beispiel: <code>#{minusYears(myDateTime,1)}</code></p> <p>Ergebnis: "2010-05-24T17:10:00z"</p>
<code>int minute(DateTime myDateTime)</code>	<p>Ruft die Minute des DateTime Werts als Ganzzahl ab.</p> <p>Beispiel: <code>#{minute(myDateTime)}</code></p> <p>Ergebnis: 10</p>

Funktion	Beschreibung
<code>int month(DateTime myDateTime)</code>	<p>Ruft den Monat des DateTime Werts als Ganzzahl ab.</p> <p>Beispiel: <code>#{month(myDateTime)}</code></p> <p>Ergebnis: 5</p>
<code>DateTime plusDays(DateTime myDateTime, int daysToAdd)</code>	<p>Erstellt ein DateTime Objekt, das sich aus der Addition der angegebenen Anzahl von Tagen zur angegebenen Anzahl von Tagen ergibt DateTime.</p> <p>Beispiel: <code>#{plusDays(myDateTime,1)}</code></p> <p>Ergebnis: "2011-05-25T17:10:00z"</p>
<code>DateTime plusHours(DateTime myDateTime, int hoursToAdd)</code>	<p>Erstellt ein DateTime Objekt, das sich aus der Addition der angegebenen Anzahl von Stunden zur angegebenen Anzahl von Stunden ergibt DateTime.</p> <p>Beispiel: <code>#{plusHours(myDateTime,1)}</code></p> <p>Ergebnis: "2011-05-24T18:10:00z"</p>

Funktion	Beschreibung
<code>DateTime plusMinutes(DateTime myDateTime,int minutesToAdd)</code>	<p>Erstellt ein DateTime Objekt, das sich aus der Addition der angegebenen Anzahl von Minuten zur angegebenen Anzahl von Minuten ergibt DateTime.</p> <p>Beispiel: <code>#{plusMinutes(myDateTime,1)}</code></p> <p>Ergebnis: "2011-05-24 17:11:00z"</p>
<code>DateTime plusMonths(DateTime myDateTime,int monthsToAdd)</code>	<p>Erstellt ein DateTime Objekt, das sich aus der Addition der angegebenen Anzahl von Monaten zur angegebenen Anzahl von Monaten ergibt DateTime.</p> <p>Beispiel: <code>#{plusMonths(myDateTime,1)}</code></p> <p>Ergebnis: "2011-06-24T17:10:00z"</p>
<code>DateTime plusWeeks(DateTime myDateTime,int weeksToAdd)</code>	<p>Erstellt ein DateTime Objekt, das das Ergebnis der Addition der angegebenen Anzahl von Wochen zur angegebenen Anzahl ist DateTime.</p> <p>Beispiel: <code>#{plusWeeks(myDateTime,1)}</code></p> <p>Ergebnis: "2011-05-31T17:10:00z"</p>

Funktion	Beschreibung
<code>DateTime plusYears(DateTime myDateTime, int yearsToAdd)</code>	<p>Erstellt ein DateTime Objekt, das sich aus der Addition der angegebenen Anzahl von Jahren zur angegebenen Anzahl von Jahren ergibt DateTime.</p> <p>Beispiel: <code>#{plusYears(myDateTime,1)}</code></p> <p>Ergebnis: "2012-05-24T17:10:00z"</p>
<code>DateTime sunday(DateTime myDateTime)</code>	<p>Erstellt ein DateTime Objekt für den vorangegangenen Sonntag, relativ zum angegebenen Wert DateTime. Wenn der angegebene Sonntag ein Sonntag DateTime ist, entspricht das Ergebnis dem angegebenen DateTime.</p> <p>Beispiel: <code>#{sunday(myDateTime)}</code></p> <p>Ergebnis: "2011-05-22 17:10:00 UTC"</p>
<code>int year(DateTime myDateTime)</code>	<p>Ruft das Jahr des DateTime Werts als Ganzzahl ab.</p> <p>Beispiel: <code>#{year(myDateTime)}</code></p> <p>Ergebnis: 2011</p>

Funktion	Beschreibung
<code>DateTime yesterday(DateTime myDateTime)</code>	<p>Erstellt ein DateTime Objekt für den vorherigen Tag, relativ zum angegebenen Wert DateTime. Das Ergebnis ist identisch mit minusDays(1).</p> <p>Beispiel: <code>#{{yesterday(myDateTime)}}</code></p> <p>Ergebnis: "2011-05-23T17:10:00z"</p>

Sonderzeichen

AWS Data Pipeline verwendet bestimmte Zeichen, die in Pipeline-Definitionen eine besondere Bedeutung haben, wie in der folgenden Tabelle gezeigt.

Sonderzeichen	Beschreibung	Beispiele
@	Laufzeitfeld. Wenn dieses Zeichen dem Namen eines Feldes vorangestellt wird, ist dieses nur während der Ausführung der Pipeline verfügbar.	<code>@actualStartTime</code> <code>@failureReason</code> <code>@resourceStatus</code>
#	Ausdruck. Ausdrücke werden durch „# {“ und „}“ getrennt, und der Inhalt der geschweiften Klammern wird durch ausgewertet. AWS Data Pipeline Weitere Informationen finden Sie unter Ausdrücke .	<code># {format (myDateTime, 'hh:mm:ss'YYYY-MM-dd)}</code> <code>s3://amzn-s3-demo-bucket#{id}.csv</code>

Sonderzeichen	Beschreibung	Beispiele
*	Verschlüsseltes Feld. Dieses Zeichen ist ein Feldnamen präfix, das angibt, dass der Inhalt dieses Felds bei der Übertragung zwischen der Konsole oder CLI und dem AWS Data Pipeline Dienst verschlüsselt werden AWS Data Pipeline soll.	*Passwort

Pipeline-Objektreferenz

Sie können die folgenden Pipeline-Objekte und -Komponenten in Ihrer Pipeline-Definitionsdatei verwenden.

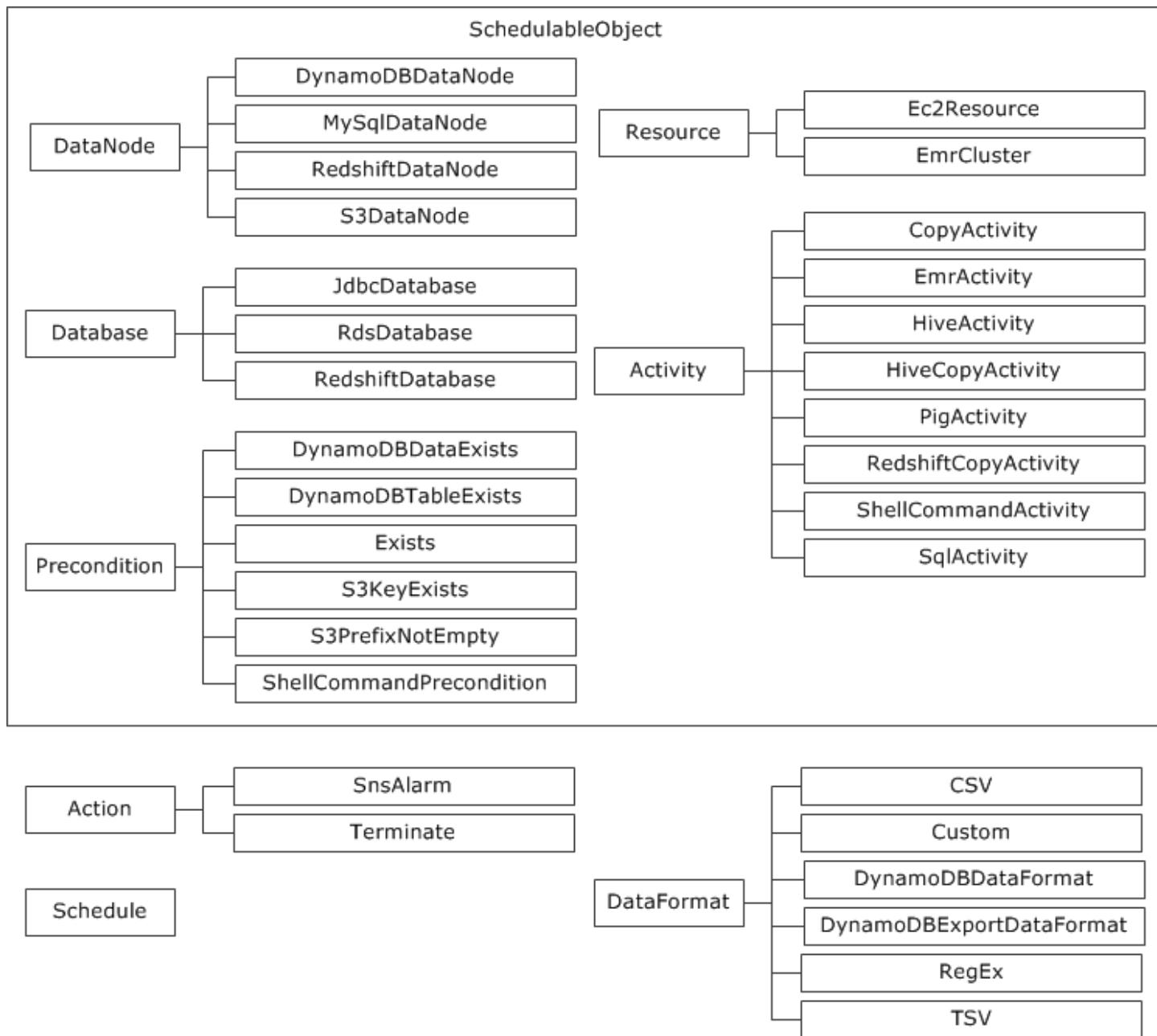
Inhalt

- [Datenknoten](#)
- [Aktivitäten](#)
- [Ressourcen](#)
- [Vorbedingungen](#)
- [Datenbanken](#)
- [Datenformate](#)
- [Aktionen](#)
- [Plan](#)
- [Dienstprogramme](#)

 Note

Eine Beispielanwendung, die das AWS Data Pipeline Java-SDK verwendet, finden Sie unter [Data Pipeline DynamoDB Export Java Sample on GitHub](#)

Im Folgenden finden Sie die Objekthierarchie für AWS Data Pipeline



Datenknoten

Im Folgenden sind die AWS Data Pipeline Datenknotenobjekte aufgeführt:

Objekte

- [Dynamo-Knoten DBData](#)
- [MySqlDataNode](#)
- [RedshiftDataNode](#)

- [S3 DataNode](#)
- [SqlDataNode](#)

Dynamo-Knoten DBData

Definiert mithilfe von DynamoDB einen Datenknoten, der als Eingabe für ein OR-Objekt angegeben wird. [HiveActivity](#) [EMRActivity](#)

 Note

Das DynamoDBDataNode-Objekt unterstützt die Vorbedingung Exists nicht.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Dieses Objekt verweist auf zwei andere Objekte, die Sie in derselben Pipeline-Definitionsdatei definieren. CopyPeriod ist ein Schedule-Objekt und Ready ist ein Vorbedingungsobjekt.

```
{
  "id" : "MyDynamoDBTable",
  "type" : "DynamoDBDataNode",
  "schedule" : { "ref" : "CopyPeriod" },
  "tableName" : "adEvents",
  "precondition" : { "ref" : "Ready" }
}
```

Syntax

Pflichtfelder	Description	Slot-Typ
tableName	Die DynamoDB-Tabelle.	String

Objektaufruf-Felder	Description	Slot-Typ
schedule	Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Benutzer	Referenzobjekt, zum Beispiel „schedule“:

Objektaufruf-Felder	Description	Slot-Typ
	<p>müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Benutzer können diese Anforderung erfüllen, indem sie explizit einen Zeitplan für das Objekt festlegen, indem sie beispielsweise „schedule“: {"ref": "DefaultSchedule,"} angeben. In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Benutzer ein übergeordnetes Objekt mit Zeitpläne referenz erstellen. Weitere Informationen zu optionalen Zeitplankonfigurationen finden Sie unter Zeitplan.</p>	{"ref": "myScheduleId" ,}

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn dieses Feld aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
dataFormat	DataFormat für die von diesem Datenknoten beschriebenen Daten. Derzeit unterstützt für HiveActivity und HiveCopyActivity.	Referenzobjekt, „DataFormat“: {"ref" DBData FormatId :"MyDynam o „}

Optionale Felder	Description	Slot-Typ
dependsOn	Angeben der Abhängigkeit von einem anderen ausführbaren Objekt	Referenzobjekt, z. B. „dependSon“: {"ref": "„} myActivityId
failureAndRerunMode	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden	Aufzählung
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlgeschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId „}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id „}
pipelineLogUri	Die S3-URI (wie 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String

Optionale Felder	Description	Slot-Typ
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. „precondition“: {"ref": "„} myPreconditionId
readThroughputPercent	Legt die Rate der Lesevorgänge so fest, dass Ihre von DynamoDB bereitgestellte Durchsatzrate im für Ihre Tabelle zugewiesenen Bereich liegt. Der Wert ist zweistellig und liegt zwischen 0,1 und 1,0 (einschließlich).	Double
Region	Der Code für die Region, in der die DynamoDB-Tabelle vorhanden ist. Beispiel: us-east-1. Dies wird verwendet, HiveActivity wenn es Staging für DynamoDB-Tabellen in Hive durchführt.	Aufzählung
reportProgressTimeout	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum
runsOn	Die Rechenressource zum Ausführen der Aktivität oder des Befehls. Zum Beispiel eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId „}

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden.</p> <p>Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.</p>	Aufzählung
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird workerGroup ignoriert.	String
writeThroughputPer cent	Legt die Rate der Schreibvorgänge so fest, dass Ihre von DynamoDB bereitgestellte Durchsatzrate im für Ihre Tabelle zugewiesenen Bereich liegt. Der Wert ist zweistellig und liegt zwischen .1 und 1.0 (einschließlich).	Double

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": " myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Zuständigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {" ref": " myRunnableObject Id "}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String

Laufzeitfelder	Description	Slot-Typ
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ healthStatusUpdated Zeit	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Laufzeitfelder	Description	Slot-Typ
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „waitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

MySqlDataNode

Legt ein Datenknoten mit MySQL fest.

 Note

Der MySqlDataNode-Typ ist veraltet. Stattdessen empfehlen wir, [SqlDataNode](#) zu verwenden.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Dieses Objekt verweist auf zwei andere Objekte, die Sie in derselben Pipeline-Definitionsdatei definieren. CopyPeriod ist ein Schedule-Objekt und Ready ist ein Vorbedingungsobjekt.

```
{
```

```

"id" : "Sql Table",
"type" : "MySQLDataNode",
"schedule" : { "ref" : "CopyPeriod" },
"table" : "adEvents",
"username": "user_name",
"*password": "my_password",
"connectionString": "jdbc:mysql://mysqlinstance-rds.example.us-east-1.rds.amazonaws.com:3306/database_name",
"selectQuery" : "select * from #{table} where eventTime >=
'#{@scheduledStartTime.format('YYYY-MM-dd HH:mm:ss')}' and eventTime <
'#{@scheduledEndTime.format('YYYY-MM-dd HH:mm:ss')}',
"precondition" : { "ref" : "Ready" }
}

```

Syntax

Pflichtfelder	Description	Slot-Typ
Tabelle	Der Name der Tabelle in der MySQL-Datenbank.	String

Objektaufruf-Felder	Description	Slot-Typ
schedule	Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Benutzer müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Benutzer können diese Anforderung erfüllen, indem sie explizit einen Zeitplan für das Objekt festlegen, indem sie beispielsweise „schedule“: {"ref": "DefaultSchedule,"} angeben. In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId"}

Objektaufruf-Felder	Description	Slot-Typ
	innerhalb des Hauptplans), können Benutzer ein übergeordnetes Objekt mit Zeitpläne referenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html .	

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
createTableSql	Ein SQL-Tabellenerstellungsausdruck, der die Tabelle erstellt.	String
Datenbank	Name der Datenbank.	Referenzobjekt, z. B. „Datenbank“: {"ref": "myDatabaseId"}
dependsOn	Gibt eine Abhängigkeit von einem anderen ausführbaren Objekt an.	Referenzobjekt, z. B. „dependSon“: {"ref": "myActivityId"}
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden.	Aufzählung

Optionale Felder	Description	Slot-Typ
insertQuery	Eine SQL-Anweisung zum Einfügen von Daten in die Tabelle.	String
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId"}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId"}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id"}
pipelineLogUri	Die S3-URI (wie 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. „precondition“: {"ref": "myPreconditionId"}

Optionale Felder	Description	Slot-Typ
reportProgressTimeout	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum
runsOn	Die Rechenressource zum Ausführen der Aktivität oder des Befehls. Zum Beispiel eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId"}

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden.</p> <p>Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.</p>	Aufzählung
schemaName	Der Name des Schemas für die Tabelle.	String
selectQuery	Eine SQL-Anweisung zum Abrufen von Daten aus der Tabelle.	String
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird workerGroup ignoriert.	String

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": " myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {" ref": " myRunnableObject Id "}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String

Laufzeitfelder	Description	Slot-Typ
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ healthStatusUpdated Zeit	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvertrag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen.	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen.	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Laufzeitfelder	Description	Slot-Typ
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „waitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Weitere Informationen finden Sie unter:

- [S3 DataNode](#)

RedshiftDataNode

Definiert einen Datenknoten mithilfe von Amazon Redshift. RedshiftDataNode stellt die Eigenschaften der Daten in einer Datenbank dar, z. B. einer Datentabelle, die von Ihrer Pipeline verwendet wird.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
  "id" : "MyRedshiftDataNode",
  "type" : "RedshiftDataNode",
```

```
"database": { "ref": "MyRedshiftDatabase" },
"tableName": "adEvents",
"schedule": { "ref": "Hour" }
}
```

Syntax

Pflichtfelder	Description	Slot-Typ
Datenbank	Die Datenbank, in der die Tabelle gespeichert ist.	Referenzobjekt, z. B. „database“: {"ref": "myRedshiftDatabaseId"}
tableName	Der Name der Amazon Redshift-Tabelle. Die Tabelle wird erstellt, falls sie noch nicht existiert und Sie sie angegeben haben createTableSql.	String

Objektaufruf-Felder	Description	Slot-Typ
schedule	Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Benutzer müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Benutzer können diese Anforderung erfüllen, indem sie explizit einen Zeitplan für das Objekt festlegen, indem sie beispielsweise „schedule“: {"ref": "DefaultSchedule"} angeben. In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Benutzer ein übergeordnetes Objekt mit Zeitpläne	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId"}

Objektaufruf-Felder	Description	Slot-Typ
	<p>nreferenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
createTableSql	<p>Ein SQL-Ausdruck, der die Tabelle in der Datenbank erstellt. Wir empfehlen, dass Sie das Schema angeben, in dem die Tabelle erstellt werden soll, zum Beispiel: CREATE TABLE mySchema.myTable (bestColumn varchar (25) primary key distkey, integer sortKey). numberOfWorks AWS Data Pipeline führt das Skript in dem createTableSql Feld aus, wenn die durch TableName angegebene Tabelle nicht in dem durch das Feld SchemaName angegebenen Schema existiert. Wenn Sie beispielsweise SchemaName als mySchema angeben, mySchema jedoch nicht in das createTableSql Feld aufnehmen, wird die Tabelle im falschen Schema erstellt (standardmäßig würde sie in PUBLIC erstellt). Dies</p>	String

Optionale Felder	Description	Slot-Typ
	passiert, da AWS Data Pipeline die CREATE-TABLE-Anweisungen nicht parst.	
dependsOn	Angeben der Abhängigkeit von einem anderen ausführbaren Objekt	Referenzobjekt, z. B. „dependSon“: {"ref": "„} myActivityId
failureAndRerunMode	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden	Aufzählung
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Die maximale Anzahl von Versuchen bei Ausfällen.	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlgeschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "„} myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "„} myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "„} myActionId „}

Optionale Felder	Description	Slot-Typ
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}
pipelineLogUri	Die S3-URI (wie 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. „precondition“: {"ref": "myPreconditionId"}
primaryKeys	Wenn Sie für eine Zieltabelle in RedShiftCopyActivity keine primaryKeys festlegen, können Sie eine Liste der Spalten angeben, die primaryKeys nutzen, die als mergeKey fungieren. Wenn Sie jedoch einen vorhandenen PrimaryKey in einer Amazon Redshift Redshift-Tabelle definiert haben, überschreibt diese Einstellung den vorhandenen Schlüssel.	String
reportProgressTime out	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum
runsOn	Die Rechenressource zum Ausführen der Aktivität oder des Befehls. Zum Beispiel eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId"}

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden.</p> <p>Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.</p>	Aufzählung
schemaName	In diesem optionalen Feld wird der Name des Schemas für die Amazon Redshift-Tabelle festgelegt. Wenn kein Name festgelegt wird, ist der Schemaname ÖFFENTLICH, was das Standardschema bei Amazon Redshift ist. Weitere Informationen finden Sie im Amazon Redshift Database Developer Guide.	String
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird workerGroup ignoriert.	String

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": " myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {" ref": " myRunnableObject Id "}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String

Laufzeitfelder	Description	Slot-Typ
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ healthStatusUpdated Zeit	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Laufzeitfelder	Description	Slot-Typ
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „waitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

S3 DataNode

Definiert einen Datenknoten mit Amazon S3. Standardmäßig DataNode verwendet der S3 serverseitige Verschlüsselung. Wenn Sie dies deaktivieren möchten, setzen Sie s3 EncryptionType auf NONE.

Note

Wenn Sie einen S3DataNode als Eingabe für CopyActivity nutzen, werden nur die Datenformate CSV und TSV unterstützt.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Dieses Objekt verweist auf ein anderes Objekt, das Sie in derselben Pipeline-Definitionsdatei definieren. CopyPeriod ist ein Schedule-Objekt.

```
{
  "id" : "OutputData",
  "type" : "S3DataNode",
  "schedule" : { "ref" : "CopyPeriod" },
  "filePath" : "s3://amzn-s3-demo-bucket/#{@scheduledStartTime}.csv"
}
```

Syntax

Objektaufruf-Felder	Description	Slot-Typ
schedule	<p>Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Benutzer müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Benutzer können diese Anforderung erfüllen, indem sie explizit einen Zeitplan für das Objekt festlegen, z. B. indem sie „schedule“: {"ref": "DefaultSchedule"} angeben. In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Benutzer ein übergeordnetes Objekt mit Zeitplänenreferenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId"}

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
Kompression	Die vom S3 beschriebene Art der Komprimierung für die DatenDataNode. „none“ ist keine Komprimierung und „gzip“ wird mit dem Gzip-Algorithmus komprimiert. Dieses Feld wird nur für die Verwendung mit Amazon Redshift und wenn Sie S3 DataNode mit CopyActivity verwenden, unterstützt.	Aufzählung
dataFormat	DataFormat für die in diesem S3 DataNode beschriebenen Daten.	Referenzobjekt, z. B. „dataFormat“: {"ref": "myDataFormat Id "}
dependsOn	Angeben der Abhängigkeit von einem anderen ausführbaren Objekt	Referenzobjekt, z. B. „dependSon“: {"ref": "myActivityId „}
directoryPath	Amazon S3 S3-Verzeichnispfad als URI: s3://my-bucket/my-key-for-directory. Sie müssen entweder einen Dateipfad (filePath) oder einen Wert für directoryPath angeben.	String
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden	Aufzählung
filePath	Der Pfad zum Objekt in Amazon S3 als URI, zum Beispiel: s3://my-bucket/my-key-for-f	String

Optionale Felder	Description	Slot-Typ
	ile. Sie müssen entweder einen Dateipfad (filePath) oder einen Wert für directoryPath angeben. Diese repräsentieren einen Ordner und einen Dateinamen. Mit dem directoryPath-Wert können Sie mehrere Dateien in einem Verzeichnis unterbringen.	
lateAfterTimeout	Die verstrichene Zeit nach dem Start der Pipeline, innerhalb derer das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
manifestFilePath	Der Amazon S3 S3-Pfad zu einer Manifestdatei in dem von Amazon Redshift unterstützten Format. AWS Data Pipeline verwendet die Manifestdatei, um die angegebenen Amazon S3 S3-Dateien in die Tabelle zu kopieren. Dieses Feld ist nur gültig, wenn a RedShiftCopyActivity auf den S3 verweistDataNode.	String
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId", }
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId", }

Optionale Felder	Description	Slot-Typ
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId „}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id „}
pipelineLogUri	Die S3-URI (wie 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. „precondition“: {"ref": „} myPreconditionId
reportProgressTime out	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum
runsOn	Die Rechenressource zum Ausführen der Aktivität oder des Befehls. Zum Beispiel eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId „}
s3 EncryptionType	Überschreibt den Amazon S3-Verschlüsselungstyp. Die Werte sind SERVER_SIDE_ENCRYPTION oder NONE. Die serverseitige Verschlüsselung ist standardmäßig aktiviert.	Aufzählung

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden.</p> <p>Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.</p>	Aufzählung
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird workerGroup ignoriert.	String

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": " myRunnableObject Id "}

Laufzeitfelder	Description	Slot-Typ
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn": {"ref": "myRunnableObject Id"}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String

Laufzeitfelder	Description	Slot-Typ
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ healthStatusUpdated Zeit	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „waitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Weitere Informationen finden Sie unter:

- [MySqlDataNode](#)

SqlDataNode

Legt ein Datenknoten mit SQL fest.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Dieses Objekt verweist auf zwei andere Objekte, die Sie in derselben Pipeline-Definitionsdatei definieren. CopyPeriod ist ein Schedule-Objekt und Ready ist ein Vorbedingungsobjekt.

```
{
  "id" : "Sql Table",
  "type" : "SqlDataNode",
  "schedule" : { "ref" : "CopyPeriod" },
  "table" : "adEvents",
  "database": "myDataBaseName",
  "selectQuery" : "select * from #{table} where eventTime >=
'#{@scheduledStartTime.format('YYYY-MM-dd HH:mm:ss')}' and eventTime <
'#{@scheduledEndTime.format('YYYY-MM-dd HH:mm:ss')}',
  "precondition" : { "ref" : "Ready" }
}
```

Syntax

Pflichtfelder	Description	Slot-Typ
Tabelle	Der Name der Tabelle in der SQL-Datenbank.	String

Objektaufruf-Felder	Description	Slot-Typ
schedule	<p>Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Benutzer müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Benutzer können diese Anforderung erfüllen, indem sie explizit einen Zeitplan für das Objekt festlegen, indem sie beispielsweise „schedule“: {"ref": "DefaultSchedule,„} angeben. In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Benutzer ein übergeordnetes Objekt mit Zeitpläne referenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId „}

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
createTableSql	Ein SQL-Tabellenerstellungsausdruck, der die Tabelle erstellt.	String
Datenbank	Name der Datenbank.	Referenzobjekt, z. B. „Datenbank“: {"ref": "myDatabaseId"}
dependsOn	Gibt die Abhängigkeit von einem anderen ausführbaren Objekt an.	Referenzobjekt, z. B. „dependSon“: {"ref": "myActivityId"}
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden.	Aufzählung
insertQuery	Eine SQL-Anweisung zum Einfügen von Daten in die Tabelle.	String
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl

Optionale Felder	Description	Slot-Typ
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId „}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id „}
pipelineLogUri	Die S3-URI (wie 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. „precondition“: {"ref": „} myPreconditionId
reportProgressTime out	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
runsOn	Die Rechenressource zum Ausführen der Aktivität oder des Befehls. Zum Beispiel eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId"}
scheduleType	Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden. Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.	Aufzählung
schemaName	Der Name des Schemas für die Tabelle.	String
selectQuery	Eine SQL-Anweisung zum Abrufen von Daten aus der Tabelle.	String

Optionale Felder	Description	Slot-Typ
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird workerGroup ignoriert.	String

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": " myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {" ref": " myRunnableObject Id "}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String

Laufzeitfelder	Description	Slot-Typ
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ healthStatusUpdated Zeit	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenversuch aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime

Laufzeitfelder	Description	Slot-Typ
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „waitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Weitere Informationen finden Sie unter:

- [S3 DataNode](#)

Aktivitäten

Im Folgenden sind die AWS Data Pipeline Aktivitätsobjekte aufgeführt:

Objekte

- [CopyActivity](#)
- [EmrActivity](#)
- [HadoopActivity](#)
- [HiveActivity](#)
- [HiveCopyActivity](#)
- [PigActivity](#)
- [RedshiftCopyActivity](#)
- [ShellCommandActivity](#)
- [SqlActivity](#)

CopyActivity

Kopiert Daten von einem Ort an einen anderen. CopyActivity unterstützt [S3 DataNode](#) und [SqlDataNode](#) als Eingabe und Ausgabe, und der Kopiervorgang wird normalerweise ausgeführt record-by-record. CopyActivity Stellt jedoch eine leistungsstarke Kopie von Amazon S3 zu Amazon S3 bereit, wenn alle folgenden Bedingungen erfüllt sind:

- Die Eingabe und Ausgabe sind S3 DataNodes
- Das Feld dataFormat ist für Ein- und Ausgabe dasselbe.

Wenn Sie die komprimierten Daten als Eingabe verwenden und dies nicht über das Feld compression auf den S3-Datenknoten angeben, kann CopyActivity möglicherweise fehlschlagen. In diesem Fall erkennt CopyActivity das Ende Datensatzzeichens nicht ordnungsgemäß und der Vorgang schlägt fehl. CopyActivity Unterstützt außerdem das Kopieren von einem Verzeichnis in ein anderes Verzeichnis und das Kopieren einer Datei in ein Verzeichnis. Das record-by-record Kopieren erfolgt jedoch, wenn ein Verzeichnis in eine Datei kopiert wird. Schließlich CopyActivity unterstützt das Kopieren mehrerer Amazon S3 S3-Dateien nicht.

Bei CopyActivity gibt es bestimmte Einschränkungen der CSV-Unterstützung. Wenn Sie ein S3 DataNode als Eingabe für verwendenCopyActivity, können Sie nur eine Unix/Linux Variante des CSV-Datendateiformats für die Amazon S3 S3-Eingabe- und Ausgabefelder verwenden. Die Unix/Linux Variante erfordert Folgendes:

- Das Trennzeichen muss ein Komma (,) sein.

- Die Datensätze werden nicht in Anführungszeichen gesetzt.
- Das Standard-Escape-Zeichen ist ASCII-Wert 92 (Backslash).
- Das Datensatzende-Identifier ist ASCII-Wert 10 (oder "\n").

Windows-basierte Systeme verwenden in der Regel eine andere end-of-record Zeichenfolge: einen Zeilenumbruch und einen Zeilenvorschub zusammen (ASCII-Wert 13 und ASCII-Wert 10). Sie müssen diesen Unterschied mit einem zusätzlichen Mechanismus ausgleichen, z. B. einem Skript zum Ändern der Eingabedaten vor dem Kopieren, um sicherzustellen, dass CopyActivity das Datensatzende korrekt erkennt. Andernfalls schlägt CopyActivity wiederholt fehl.

Wenn Sie mit CopyActivity einen Exportvorgang von einem PostgreSQL-RDS-Objekt in das TSV-Datenformat durchführen, ist das Standard-NUL-Zeichen \n.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Dieses Objekt verweist auf drei andere Objekte, die Sie in derselben Pipeline-Definitionsdatei definieren. CopyPeriod ist ein Schedule-Objekt und InputData und OutputData sind Datenknotenobjekte.

```
{
  "id" : "S3ToS3Copy",
  "type" : "CopyActivity",
  "schedule" : { "ref" : "CopyPeriod" },
  "input" : { "ref" : "InputData" },
  "output" : { "ref" : "OutputData" },
  "runsOn" : { "ref" : "MyEc2Resource" }
}
```

Syntax

Objektaufruf-Felder	Description	Slot-Typ
schedule	Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Benutzer müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId"}

Objektaufruf-Felder	Description	Slot-Typ
	<p>en. Benutzer können diese Anforderung erfüllen, indem sie explizit einen Zeitplan für das Objekt festlegen, indem sie beispielsweise „schedule“: {"ref": ""} angeben. DefaultSchedule In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Benutzer ein übergeordnetes Objekt mit Zeitpläne referenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
runsOn	Die Rechenressource zum Ausführen der Aktivität oder des Befehls. Zum Beispiel eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId"}
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird workerGroup ignoriert.	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
dependsOn	Angeben der Abhängigkeit von einem anderen ausführbaren Objekt.	Referenzobjekt, z. B. „dependSon“: {"ref": "myActivityId"}
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden	Aufzählung
input	Die Eingangsdatenquelle.	Referenzobjekt, z. B. „input“: {"ref": "myDataNode Id"}
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId"}

Optionale Felder	Description	Slot-Typ
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. "onSuccess": {"ref": "myActionId"}
output	Die Eingangsdatenquelle.	Referenzobjekt, z. B. "output": {"ref": "myDataNode Id"}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. "parent": {"ref": "myBaseObject Id"}
pipelineLogUri	Die S3-URI (wie 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. "precondition": {"ref": "myPreconditionId"}
reportProgressTime out	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden.</p> <p>Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.</p>	Aufzählung

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": " myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime

Laufzeitfelder	Description	Slot-Typ
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Zuständigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {" ref": " myRunnableObject Id "}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String

Laufzeitfelder	Description	Slot-Typ
@ healthStatusUpdated Zeit	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „waitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Weitere Informationen finden Sie unter:

- [ShellCommandActivity](#)
- [EmrActivity](#)
- [Exportieren von MySQL-Daten nach Amazon S3 mithilfe AWS Data Pipeline](#)

EmrActivity

Führt einen EMR-Cluster.

AWS Data Pipeline verwendet ein anderes Format für Schritte als Amazon EMR; AWS Data Pipeline verwendet beispielsweise kommagetrennte Argumente nach dem JAR-Namen im `EmrActivity` Schrittfeld. Das folgende Beispiel zeigt einen für Amazon EMR formatierten Schritt, gefolgt von seinem AWS Data Pipeline Äquivalent:

```
s3://amzn-s3-demo-bucket/MyWork.jar arg1 arg2 arg3
```

```
"s3://amzn-s3-demo-bucket/MyWork.jar,arg1,arg2,arg3"
```

Beispiele

Es folgt ein Beispiel für diesen Objekttyp. In diesem Beispiel werden ältere Versionen von Amazon EMR verwendet. Überprüfen Sie die Richtigkeit dieses Beispiels anhand der Version des Amazon EMR-Clusters, die Sie verwenden.

Dieses Objekt verweist auf drei andere Objekte, die Sie in derselben Pipeline-Definitionsdatei definieren. MyEmrCluster ist ein EmrCluster-Objekt und MyS3Input und MyS3Output sind S3DataNode-Objekte.

Note

In diesem Beispiel können Sie das Feld step mit der gewünschten Cluster-Zeichenfolge ersetzen. Hierbei kann es sich u. a. um ein Pig-Skript, ein Hadoop-Streaming-Cluster oder Ihre eigene benutzerdefinierte JAR-Datei mit ihren Parametern handeln.

Hadoop 2.x (AMI 3.x)

```
{  
  "id" : "MyEmrActivity",  
  "type" : "EmrActivity",  
  "runsOn" : { "ref" : "MyEmrCluster" },  
  "preStepCommand" : "scp remoteFiles localFiles",  
  "step" : ["s3://amzn-s3-demo-bucket/myPath/myStep.jar,firstArg,secondArg,-files,s3://  
amzn-s3-demo-bucket/myPath/myFile.py,-input,s3://myinputbucket/path,-output,s3://  
myoutputbucket/path,-mapper,myFile.py,-reducer,reducerName","s3://amzn-s3-demo-bucket/  
myPath/myotherStep.jar,..."],  
  "postStepCommand" : "scp localFiles remoteFiles",  
  "input" : { "ref" : "MyS3Input" },  
  "output" : { "ref" : "MyS3Output" }  
}
```

Note

Um in einem Schritt Argumente an eine Anwendung zu übergeben, müssen Sie die Region im Pfad des Skripts angeben, wie im folgenden Beispiel gezeigt: Darüber hinaus müssen Sie für die zu übergebenden Argumente möglicherweise ein Escape-Zeichen verwenden. Wenn Sie beispielsweise mit `script-runner.jar` ein Shell-Skript ausführen und Argumente an das Skript übergeben möchten, müssen Sie für die Kommas, die als Trennzeichen dienen, Escape-Zeichen verwenden. Der folgende Schritt-Slot veranschaulicht die entsprechende Vorgehensweise:

```
"step" : "s3://eu-west-1.elasticmapreduce/libs/script-runner/script-  
runner.jar,s3://datapipeline/echo.sh,a\\\",b\\\\\\\",c"
```

Dieser Schritt verwendet `script-runner.jar`, um das Shell-Skript `echo.sh` auszuführen, und übergibt a, b und c als einzelne Argumente an das Skript. Die erste Escape-Zeichen wird vom resultierenden Argument entfernt, weshalb Sie möglicherweise erneut ein Escape-Zeichen verwenden müssen. Wenn Sie beispielsweise `File\\.gz` als Argument in JSON verwendet haben, können Sie als Escape-Zeichen `File\\\\.gz` verwenden. Da das erste Escape-Zeichen jedoch verworfen wird, müssen Sie `File\\\\\\\\\\\\.gz` verwenden.

Syntax

Objektaufruf-Felder	Description	Slot-Typ
schedule	<p>Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Sie müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihefolge für dieses Objekt festzulegen. Sie können diese Anforderung erfüllen, indem Sie explizit einen Zeitplan für das Objekt festlegen, indem sie beispielsweise <code>"schedule": {"ref": "DefaultSchedule"}</code> angeben. In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Sie ein übergeordnetes Objekt mit Zeitplänenreferenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	Referenzobjekt, zum Beispiel „schedule“: <code>{"ref": "myScheduleId"}</code>

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
runsOn	Der Amazon EMR-Cluster, auf dem dieser Job ausgeführt wird.	Referenzobjekt, zum Beispiel „runsOn“: {"ref": " myEmrClusterId "}
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird ignoriert.workerGroup	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
dependsOn	Angeben der Abhängigkeit von einem anderen ausführbaren Objekt.	Referenzobjekt, zum Beispiel „dependSon“: {"ref": " myActivityId „}
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden.	Aufzählung
input	Der Speicherort der Eingabedaten.	Referenzobjekt, zum Beispiel „input“:

Optionale Felder	Description	Slot-Typ
		{"ref": " myDataNode Id "}
lateAfterTimeout	Die verstrichene Zeit nach dem Start der Pipeline, innerhalb derer das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Die maximale Anzahl von Versuchen bei Ausfällen.	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, zum Beispiel „onFail“: {"ref": " myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, zum Beispiel "onLateAction": {"ref": " myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, zum Beispiel „onSuccess“: {"ref": " myActionId „}
output	Der Speicherort der Ausgabedaten.	Referenzobjekt, zum Beispiel „output“: {"ref": " myDataNode Id "}

Optionale Felder	Description	Slot-Typ
übergeordneter	Das übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, zum Beispiel „parent“: {"ref": " myBaseObject Id "}
pipelineLogUri	Die Amazon S3 S3-URI, z. B. 's3://BucketName/Prefix/ 'für das Hochladen von Protokollen für die Pipeline.	String
postStepCommand	Shell-Skripts, die nach Abschluss aller Schritte ausgeführt werden. Wenn Sie mehrere Skripts angeben möchten (maximal 255), fügen Sie die entsprechende Anzahl von postStepCommand -Feldern hinzu.	String
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, zum Beispiel „precondition“: {"ref": " myPreconditionId "}
preStepCommand	Shell-Skripts, die vor allen Schritten ausgeführt werden. Wenn Sie mehrere Skripts angeben möchten (maximal 255), fügen Sie die entsprechende Anzahl von preStepCommand -Feldern hinzu.	String
reportProgressTimeout	Das Timeout für aufeinanderfolgende Aufrufe von reportProgress durch Remote-Arbeit. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum

Optionale Felder	Description	Slot-Typ
resizeClusterBeforeRunning	Ändern Sie die Größe des Clusters, bevor Sie diese Aktivität ausführen, um DynamoDB-Tabellen aufzunehmen, die als Eingaben oder Ausgaben angegeben sind.	Boolesch
	<p> Note</p> <p>Wenn Sie <code>EmrActivity</code> entweder <code>DynamoDBDataNode</code> als Eingabe- oder Ausgabeknoten verwenden und den Wert auf <code>setzenTRUE</code>, AWS Data Pipeline beginnen Sie mit der <code>resizeClusterBeforeRunning</code> Verwendung von <code>m3.xlarge</code> Instance-Typen. Dadurch wird Ihre Auswahl an Instance-Typen mit <code>m3.xlarge</code> überschrieben, wodurch Ihre monatlichen Kosten ansteigen könnten.</p>	
resizeClusterMaxInstances	Ein Limit für die maximale Anzahl von Instances, die vom Resize-Algorithmus angefordert werden können.	Ganzzahl
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Werte sind <code>cron</code>, <code>ondemand</code> und <code>timeseries</code>. Die <code>timeseries</code>-Planung bedeutet, dass Instances am Ende jedes Intervalls geplant sind. Die <code>cron</code>-Planung bedeutet, dass Instances am Anfang jedes Intervalls geplant sind. Ein <code>ondemand</code>-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Sie müssen die Pipeline nicht klonen oder neu erstellen, um sie erneut auszuführen. Wenn Sie einen <code>ondemand</code>-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene <code>scheduleType</code> sein. Um <code>ondemand</code>-Pipelines zu verwenden, rufen Sie einfach den <code>ActivatePipeline</code>-Vorgang für jeden nachfolgenden Lauf auf.</p>	Aufzählung
Schritt	<p>Einzelne oder mehrere vom Cluster auszuführende Schritte. Wenn Sie mehrere Schritte angeben möchten (maximal 255), fügen Sie die entsprechende Anzahl von <code>step</code>-Feldern hinzu. Verwenden Sie durch Komma getrennte Argumente nach dem JAR-Namen, z. B. "s3://amzn-s3-demo-bucket/MyWork.jar,arg1,arg2,arg3".</p>	String

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „activeInstances“: {"ref": " myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, zum Beispiel "cascadeFailedOn": {" ref": " myRunnableObject Id "}
emrStepLog	Amazon EMR-Schrittprotokolle sind nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage , wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String

Laufzeitfelder	Description	Slot-Typ
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ Zeit healthStatusUpdated	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvertrag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für das Objekt einplanen.	DateTime
@scheduledStartTime	Startzeit für das Objekt einplanen.	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde	String

Laufzeitfelder	Description	Slot-Typ
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, zum Beispiel „WaitingOn“: {"ref": "myRunnableObject Id "}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Weitere Informationen finden Sie unter:

- [ShellCommandActivity](#)
- [CopyActivity](#)
- [EmrCluster](#)

HadoopActivity

Führt einen MapReduce Job auf einem Cluster aus. Bei dem Cluster kann es sich um einen EMR-Cluster handeln, der von AWS Data Pipeline oder einer anderen Ressource verwaltet wird, wenn Sie ihn verwenden TaskRunner. Verwenden Sie diese Option, HadoopActivity wenn Sie parallel arbeiten möchten. Auf diese Weise können Sie die Planungsressourcen des YARN-Frameworks oder des MapReduce Resource Negotiators in Hadoop 1 verwenden. Wenn Sie die Arbeit sequenziell mit der Amazon EMR Step-Aktion ausführen möchten, können Sie dies trotzdem verwenden. [EmrActivity](#)

Beispiele

HadoopActivity unter Verwendung eines EMR-Clusters, verwaltet von AWS Data Pipeline

Das folgende HadoopActivity Objekt verwendet eine EmrCluster Ressource, um ein Programm auszuführen:

```
{
  "name": "MyHadoopActivity",
  "schedule": {"ref": "ResourcePeriod"},
  "runsOn": {"ref": "MyEmrCluster"},
  "type": "HadoopActivity",
  "preActivityTaskConfig": {"ref": "preTaskScriptConfig"},
  "jarUri": "/home/hadoop/contrib/streaming/hadoop-streaming.jar",
  "argument": [
    "-files",
    "s3://elasticmapreduce/samples/wordcount/wordSplitter.py",
    "-mapper",
    "wordSplitter.py",
    "-reducer",
    "aggregate",
    "-input",
    "s3://elasticmapreduce/samples/wordcount/input/",
    "-output",
    "s3://amzn-s3-demo-bucket/MyHadoopActivity/#{@pipelineId}/
#{format(@scheduledStartTime, 'YYYY-MM-dd')}"
  ],
  "maximumRetries": "0",
  "postActivityTaskConfig": {"ref": "postTaskScriptConfig"},
  "hadoopQueue" : "high"
}
```

Hier ist das entsprechende *MyEmrCluster*, das die FairScheduler und -Warteschlangen in YARN für Hadoop 2-basiert konfiguriert: AMIs

```
{
  "id" : "MyEmrCluster",
  "type" : "EmrCluster",
  "hadoopSchedulerType" : "PARALLEL_FAIR_SCHEDULING",
  "amiVersion" : "3.7.0",
  "bootstrapAction" : ["s3://Region.elasticmapreduce/bootstrap-
actions/configure-hadoop,-z,yarn.scheduler.capacity.root.queues=low
\,high\,default,-z,yarn.scheduler.capacity.root.high.capacity=50,-
```

```

z,yarn.scheduler.capacity.root.low.capacity=10,-
z,yarn.scheduler.capacity.root.default.capacity=30"]
}

```

Dies ist der, den EmrCluster Sie zur Konfiguration FairScheduler in Hadoop 1 verwenden:

```

{
  "id": "MyEmrCluster",
  "type": "EmrCluster",
  "hadoopSchedulerType": "PARALLEL_FAIR_SCHEDULING",
  "amiVersion": "2.4.8",
  "bootstrapAction": "s3://Region.elasticmapreduce/bootstrap-
actions/configure-hadoop,-m,mapred.queue.names=low\\\\\\,high\\\\\\,default,-
m,mapred.fairscheduler.poolnameproperty=mapred.job.queue.name"
}

```

Die folgenden Konfigurationen EmrCluster basieren auf CapacityScheduler Hadoop 2: AMIs

```

{
  "id": "MyEmrCluster",
  "type": "EmrCluster",
  "hadoopSchedulerType": "PARALLEL_CAPACITY_SCHEDULING",
  "amiVersion": "3.7.0",
  "bootstrapAction": "s3://Region.elasticmapreduce/bootstrap-
actions/configure-hadoop,-z,yarn.scheduler.capacity.root.queues=low
\\\\\\,high,-z,yarn.scheduler.capacity.root.high.capacity=40,-
z,yarn.scheduler.capacity.root.low.capacity=60"
}

```

HadoopActivity mit einem vorhandenen EMR-Cluster

In diesem Beispiel verwenden Sie workergroups und a, TaskRunner um ein Programm auf einem vorhandenen EMR-Cluster auszuführen. Die folgende Pipeline-Definition dient dazu: HadoopActivity

- Führen Sie ein MapReduce Programm nur auf *myWorkerGroup* Ressourcen aus. Weitere Informationen zu Worker-Gruppen finden Sie unter [Arbeiten an vorhandenen Ressourcen mithilfe von Task Runner ausführen](#).
- Führen Sie eine preActivityTask Config und eine postActivityTask Config aus

```
{
}
```

```
"objects": [
  {
    "argument": [
      "-files",
      "s3://elasticmapreduce/samples/wordcount/wordSplitter.py",
      "-mapper",
      "wordSplitter.py",
      "-reducer",
      "aggregate",
      "-input",
      "s3://elasticmapreduce/samples/wordcount/input/",
      "-output",
      "s3://amzn-s3-demo-bucket/MyHadoopActivity/#{@pipelineId}/
#{format(@scheduledStartTime, 'YYYY-MM-dd')}"
    ],
    "id": "MyHadoopActivity",
    "jarUri": "/home/hadoop/contrib/streaming/hadoop-streaming.jar",
    "name": "MyHadoopActivity",
    "type": "HadoopActivity"
  },
  {
    "id": "SchedulePeriod",
    "startDateTime": "start_datetime",
    "name": "SchedulePeriod",
    "period": "1 day",
    "type": "Schedule",
    "endDateTime": "end_datetime"
  },
  {
    "id": "ShellScriptConfig",
    "scriptUri": "s3://amzn-s3-demo-bucket/scripts/preTaskScript.sh",
    "name": "preTaskScriptConfig",
    "scriptArgument": [
      "test",
      "argument"
    ],
    "type": "ShellScriptConfig"
  },
  {
    "id": "ShellScriptConfig",
    "scriptUri": "s3://amzn-s3-demo-bucket/scripts/postTaskScript.sh",
    "name": "postTaskScriptConfig",
    "scriptArgument": [
      "test",
      "argument"
    ],
    "type": "ShellScriptConfig"
  }
]
```

```

    "argument"
  ],
  "type": "ShellScriptConfig"
},
{
  "id": "Default",
  "scheduleType": "cron",
  "schedule": {
    "ref": "SchedulePeriod"
  },
  "name": "Default",
  "pipelineLogUri": "s3://amzn-s3-demo-bucket/
logs/2015-05-22T18:02:00.343Z642f3fe415",
  "maximumRetries": "0",
  "workerGroup": "myWorkerGroup",
  "preActivityTaskConfig": {
    "ref": "preTaskScriptConfig"
  },
  "postActivityTaskConfig": {
    "ref": "postTaskScriptConfig"
  }
}
]
}

```

Syntax

Pflichtfelder	Description	Slot-Typ
jarUri	Speicherort einer JAR in Amazon S3 oder im lokalen Dateisystem des Clusters, mit dem ausgeführt werden soll HadoopActivity.	String

Objektaufruf-Felder	Description	Slot-Typ
schedule	Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Benutzer müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausfü	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId"}

Objektaufruf-Felder	Description	Slot-Typ
	<p>hrungsreihenfolge für dieses Objekt festzulegen. Benutzer können diese Anforderung erfüllen, indem sie explizit einen Zeitplan für das Objekt festlegen, z. B. indem sie „schedule“: {"ref": "DefaultSchedule,"} angeben. In den meisten Fällen ist es besser, den Zeitplan zuweisen auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Benutzer ein übergeordnetes Objekt mit Zeitplänenreferenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
runsOn	EMR-Cluster, auf dem dieser Auftrag ausgeführt wird.	Referenzobjekt, z. B. „runsOn“: {"ref": "myEmrCluster Id "}
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird workerGroup ignoriert.	String

Optionale Felder	Description	Slot-Typ
argument	Argumente, die an die JAR-Dateien übergeben werden.	String
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
dependsOn	Angeben der Abhängigkeit von einem anderen ausführbaren Objekt.	Referenzobjekt, z. B. „dependSon“: {"ref": "myActivityId"}
failureAndRerunodus	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden	Aufzählung
hadoopQueue	Der Name der Hadoop-Scheduler-Warteschlange, an die die Aktivität übergeben wird.	String
input	Speicherort der Eingabedaten.	Referenzobjekt, z. B. „input“: {"ref": "myDataNode Id"}
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf ondemand eingestellt ist.	Zeitraum
mainClass	Die Hauptklasse der JAR, mit der Sie die Ausführung ausführen HadoopActivity.	String

Optionale Felder	Description	Slot-Typ
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId „}
output	Speicherort der Ausgabedaten.	Referenzobjekt, z. B. „output“: {"ref": "myDataNode Id "}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}
pipelineLogUri	Die S3-URI (wie 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String
postActivityTaskConfig	Post-Activity-Konfigurationsskript, das ausgeführt werden soll. Dieses besteht aus einer URI des Shell-Skripts in Amazon S3 und einer Liste von Argumenten.	Referenzobjekt, z. B. "postActivityTaskConfig": {"ref": "myShellScriptConfigId „}

Optionale Felder	Description	Slot-Typ
preActivityTaskConfig	Pre-Activity-Konfigurationsskript, das ausgeführt werden soll. Dieses besteht aus einer URI des Shell-Skripts in Amazon S3 und einer Liste von Argumenten.	Referenzobjekt, z. B. "preActivityTaskConfig": {"ref": " myShellScriptConfigId „}
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. „Vorbedingung“: {"ref": " myPreconditionId „}
reportProgressTimeout	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden.</p> <p>Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.</p>	Aufzählung

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": " myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime

Laufzeitfelder	Description	Slot-Typ
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {" ref": " myRunnableObject Id "}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String

Laufzeitfelder	Description	Slot-Typ
@ healthStatusUpdated Zeit	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „waitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Weitere Informationen finden Sie unter:

- [ShellCommandActivity](#)
- [CopyActivity](#)
- [EmrCluster](#)

HiveActivity

Führt eine Hive-Abfrage auf einem EMR-Cluster aus. HiveActivity erleichtert die Einrichtung einer Amazon EMR-Aktivität und erstellt automatisch Hive-Tabellen auf der Grundlage von Eingabedaten, die entweder von Amazon S3 oder Amazon RDS stammen. Sie müssen lediglich den HiveQL angeben, der auf den Quelldaten ausgeführt werden soll. AWS Data Pipeline erstellt automatisch Hive-Tabellen mit \${input1}\${input2}, usw. auf der Grundlage der Eingabefelder im Objekt. **HiveActivity**

Für Amazon S3 S3-Eingaben wird das `dataFormat` Feld verwendet, um die Hive-Spaltennamen zu erstellen.

Bei MySQL-Eingaben (Amazon RDS) werden die Spaltennamen für die SQL-Abfrage verwendet, um die Hive-Spaltennamen zu erstellen.

 Note

Diese Aktivität verwendet den [CSV-Serde](#) von Hive.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Dieses Objekt verweist auf drei andere Objekte, die Sie in derselben Pipeline-Definitionsdatei definieren. MySchedule ist ein Schedule-Objekt und MyS3Input und MyS3Output sind Datenknotenobjekte.

```
{
  "name" : "ProcessLogData",
  "id" : "MyHiveActivity",
  "type" : "HiveActivity",
  "schedule" : { "ref": "MySchedule" },
  "hiveScript" : "INSERT OVERWRITE TABLE ${output1} select
host,user,time,request,status,size from ${input1};",
  "input" : { "ref": "MyS3Input" },
  "output" : { "ref": "MyS3Output" },
  "runsOn" : { "ref": "MyEmrCluster" }
}
```

Syntax

Objektaufruf-Felder	Description	Slot-Typ
schedule	Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Sie müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Sie können diese Anforderung erfüllen, indem Sie explizit einen Zeitplan für das Objekt festlegen, indem Sie beispielsweise „schedule“: {"ref": "DefaultSchedule"} angeben. In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Sie ein übergeordnetes Objekt mit Zeitplänenreferenz erstellen. Weitere Informationen zu optionalen Beispiel-	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId" „}

Objektaufruf-Felder	Description	Slot-Typ
	Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html .	

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
hiveScript	Das auszuführende Hive-Skript.	String
scriptUri	Der Speicherort des auszuführenden Hive-Skripts (z. B. s3://scriptLocation).	String

Erforderliche Gruppe	Description	Slot-Typ
runsOn	Der EMR-Cluster, auf dem diese HiveActivity ausgeführt wird	Referenzobjekt, z. B. „runsOn“: {"ref": "myEmrCluster Id "}
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird ignoriert.workerGroup	String
input	Die Eingangsdatenquelle.	Referenzobjekt, z. B. „input“: {"ref": "myDataNode Id "}
output	Die Eingangsdatenquelle.	Referenzobjekt, z. B. „output“: {"ref": "myDataNode Id "}

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
dependsOn	Angeben der Abhängigkeit von einem anderen ausführbaren Objekt.	Referenzobjekt, z. B. „dependSon“: {"ref": "myActivityId"}
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden.	Aufzählung
hadoopQueue	Der Name der Hadoop-Scheduler-Warteschlange, in der der Auftrag übermittelt wird.	String
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Die maximale Anzahl von Versuchen bei Ausfällen.	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlgeschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId"}

Optionale Felder	Description	Slot-Typ
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. "onSuccess": {"ref": "myActionId"}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. "parent": {"ref": "myBaseObject Id"}
pipelineLogUri	Die S3-URI (z. B. 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String
postActivityTaskConfig	Post-Activity-Konfigurationsskript, das ausgeführt werden soll. Dieses besteht aus einer URI des Shell-Skripts in Amazon S3 und einer Liste von Argumenten.	Referenzobjekt, z. B. "postActivityTaskConfig": {"ref": "myShellScriptConfigId"}
preActivityTaskConfig	Pre-Activity-Konfigurationsskript, das ausgeführt werden soll. Dieses besteht aus einer URI des Shell-Skripts in Amazon S3 und einer Liste von Argumenten.	Referenzobjekt, z. B. "preActivityTaskConfig": {"ref": "myShellScriptConfigId"}
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. "Vorbedingung": {"ref": "myPreconditionId"}

Optionale Felder	Description	Slot-Typ
reportProgressTimeout	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in <code>reportProgress</code> . Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
resizeClusterBeforeRunning	Ändern Sie die Größe des Clusters, bevor Sie diese Aktivität ausführen, um DynamoDB-Datenknoten aufzunehmen, die als Eingaben oder Ausgaben angegeben sind.	Boolesch
	<p> Note</p> <p>Wenn Ihre Aktivität <code>a</code> entweder <code>DynamoDBDataNode</code> als Eingabe- oder Ausgabeknoten verwendet und Sie den Wert auf <code>setzenTRUE</code>, AWS Data Pipeline beginnt die <code>resizeClusterBeforeRunning</code> Verwendung von <code>m3.xlarge</code> Instance-Typen. Dadurch wird Ihre Auswahl an Instance-Typen mit <code>m3.xlarge</code> überschrieben, wodurch Ihre monatlichen Kosten ansteigen könnten.</p>	
resizeClusterMaxInstances	Ein Limit für die maximale Anzahl von Instances, die vom Resize-Algorithmus angefordert werden können.	Ganzzahl
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden.</p> <p>Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.</p>	Aufzählung

Optionale Felder	Description	Slot-Typ
scriptVariable	Gibt Skriptvariablen an, die Amazon EMR bei der Ausführung eines Skripts an Hive weitergibt. Im folgenden Beispiel etwa würden Skriptvariablen eine SAMPLE- und FILTER_DATE-Variable an Hive übergeben: SAMPLE=s3://elasticmapreduce/samples/hive-ads und FILTER_DATE=#{format(@scheduledStartTime, 'YYYY-MM-dd')}%. Dieses Feld akzeptiert mehrere Werte und funktioniert sowohl mit script- als auch mit scriptUri -Feldern. Darüber hinaus funktioniert scriptVariable unabhängig davon, ob "stage" auf true oder false festgelegt ist. Dieses Feld ist besonders nützlich, um mithilfe von AWS Data Pipeline -Ausdrücken und -Funktionen dynamische Werte an Hive zu senden.	String
stage	Legt fest, ob vor oder nach dem Ausführen des Skripts Staging aktiviert wird. Ist mit Hive 11 unzulässig. Verwenden Sie daher eine Amazon EMR-AMI in der Version 3.2.0 oder höher.	Boolesch

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": "Id"} myRunnableObject
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime

Laufzeitfelder	Description	Slot-Typ
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {" ref": " myRunnableObject Id "}
emrStepLog	Amazon EMR-Schrittprotokolle sind nur bei EMR-Aktivitätsversuchen verfügbar.	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String

Laufzeitfelder	Description	Slot-Typ
@ Zeit healthStatusUpdated	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für ein Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für ein Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „WaitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Weitere Informationen finden Sie unter:

- [ShellCommandActivity](#)
- [EmrActivity](#)

HiveCopyActivity

Führt eine Hive-Abfrage auf einem EMR-Cluster aus. HiveCopyActivity erleichtert das Kopieren von Daten zwischen DynamoDB-Tabellen. HiveCopyActivity akzeptiert eine HiveQL-Anweisung zum Filtern von Eingabedaten aus DynamoDB auf Spalten- und Zeilenebene.

Beispiel

Das folgende Beispiel zeigt, wie Sie mit HiveCopyActivity und DynamoDBExportDataFormat Daten von einem DynamoDBDataNode auf einen anderen kopieren können, während gleichzeitig Daten basierend auf einem Zeitstempel gefiltert werden.

```
{
  "objects": [
    {
      "id" : "DataFormat.1",
      "name" : "DataFormat.1",
      "type" : "DynamoDBExportDataFormat",
      "column" : "timeStamp BIGINT"
    },
    {
      "id" : "DataFormat.2",
      "name" : "DataFormat.2",
      "type" : "DynamoDBExportDataFormat",
      "column" : "timeStamp BIGINT"
    }
  ]
}
```

```
"id" : "DataFormat.2",
"name" : "DataFormat.2",
"type" : "DynamoDBExportDataFormat"
},
{
  "id" : "DynamoDBDataNode.1",
  "name" : "DynamoDBDataNode.1",
  "type" : "DynamoDBDataNode",
  "tableName" : "item_mapped_table_restore_temp",
  "schedule" : { "ref" : "ResourcePeriod" },
  "dataFormat" : { "ref" : "DataFormat.1" }
},
{
  "id" : "DynamoDBDataNode.2",
  "name" : "DynamoDBDataNode.2",
  "type" : "DynamoDBDataNode",
  "tableName" : "restore_table",
  "region" : "us_west_1",
  "schedule" : { "ref" : "ResourcePeriod" },
  "dataFormat" : { "ref" : "DataFormat.2" }
},
{
  "id" : "EmrCluster.1",
  "name" : "EmrCluster.1",
  "type" : "EmrCluster",
  "schedule" : { "ref" : "ResourcePeriod" },
  "masterInstanceType" : "m1.xlarge",
  "coreInstanceCount" : "4"
},
{
  "id" : "HiveTransform.1",
  "name" : "Hive Copy Transform.1",
  "type" : "HiveCopyActivity",
  "input" : { "ref" : "DynamoDBDataNode.1" },
  "output" : { "ref" : "DynamoDBDataNode.2" },
  "schedule" : { "ref" : "ResourcePeriod" },
  "runsOn" : { "ref" : "EmrCluster.1" },
  "filterSql" : "`timeStamp` > unix_timestamp(\"#{@scheduledStartTime}\", \"yyyy-MM-dd'T'HH:mm:ss\")"
},
{
  "id" : "ResourcePeriod",
  "name" : "ResourcePeriod",
  "type" : "Schedule",
```

```

    "period" : "1 Hour",
    "startDateTime" : "2013-06-04T00:00:00",
    "endDateTime" : "2013-06-04T01:00:00"
}
]
}

```

Syntax

Objektaufruf-Felder	Description	Slot-Typ
schedule	<p>Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Benutzer müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Benutzer können diese Anforderung erfüllen, indem sie explizit einen Zeitplan für das Objekt festlegen, indem sie beispielsweise „schedule“: {"ref": „} angeben. DefaultSchedule In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Benutzer ein übergeordnetes Objekt mit Zeitpläne referenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId „}

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
runsOn	Geben Sie den Cluster an, auf dem ausgeführt werden soll.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId"}
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird ignoriert.	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Der zuletzt gemeldete Status von der Remote-Aktivität.	String
attemptTimeout	Das Timeout für die Fertigstellung der Remote-Arbeit. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
dependsOn	Gibt die Abhängigkeit von einem anderen ausführbaren Objekt an.	Referenzobjekt, z. B. „dependSon“: {"ref": "myActivityId"}
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden.	Aufzählung
filterSql	Ein Hive-SQL-Anweisungsfragment, das eine Teilmenge der zu kopierenden DynamoDB- oder Amazon S3 S3-Daten filtert. Der Filter	String

Optionale Felder	Description	Slot-Typ
	sollte nur Prädikate enthalten und nicht mit einer WHERE Klausel beginnen, da er automatisch hinzugefügt wird. AWS Data Pipeline	
input	Die Eingangsdatenquelle. Dies muss ein S3DataNode oder DynamoDBDataNode sein. Wenn Sie DynamoDBNode verwenden, geben Sie ein DynamoDBExportData Format an.	Referenzobjekt, z. B. „input“: {"ref": "myDataNode Id "}
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Die maximale Anzahl von Versuchen bei Ausfällen.	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction,: {"ref": "myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId „}

Optionale Felder	Description	Slot-Typ
output	Die Eingangsdatenquelle. Wenn die Eingabe S3DataNode ist, muss diese auf DynamoDBDataNode festgelegt sein. Andernfalls kann dies S3DataNode oder DynamoDBDataNode sein. Wenn Sie DynamoDBNode verwenden, geben Sie ein DynamoDBExportDataFormat an.	Referenzobjekt, z. B. „output“: {"ref": "myDataNode Id "}
übergeordneter	Das übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}
pipelineLogUri	Die Amazon S3 S3-URI, z. B. 's3://BucketName/Key/' für das Hochladen von Protokollen für die Pipeline.	String
postActivityTaskConfig	Das Post-Activity-Konfigurationsskript, das ausgeführt werden soll. Dieses besteht aus einer URI des Shell-Skripts in Amazon S3 und einer Liste von Argumenten.	Referenzobjekt, z. B. "postActivityTaskConfig": {"ref": "myShellScriptConfigId"}
preActivityTaskConfig	Das Pre-Activity-Konfigurationsskript, das ausgeführt werden soll. Dieses besteht aus einer URI des Shell-Skripts in Amazon S3 und einer Liste von Argumenten.	Referenzobjekt, z. B. "preActivityTaskConfig": {"ref": "myShellScriptConfigId"}
precondition	Definiert optional eine Vorbedingung. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. „Vorbedingung“: {"ref": "myPreconditionId"}

Optionale Felder	Description	Slot-Typ
reportProgressTimeout	Das Timeout für aufeinanderfolgende Aufrufe von <code>reportProgress</code> durch Remote-Arbeit. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
resizeClusterBeforeItIsRun	Ändern Sie die Größe des Clusters, bevor Sie diese Aktivität ausführen, um DynamoDB-Datenknoten aufzunehmen, die als Eingaben oder Ausgaben angegeben sind.	Boolesch
<p>Note</p> <p>Wenn Ihre Aktivität <code>a</code> entweder <code>DynamoDBDataNode</code> als Eingabe- oder Ausgabeknoten verwendet und Sie den Wert auf <code>setzenTRUE</code>, AWS Data Pipeline beginnt die <code>resizeClusterBeforeRunning</code> Verwendung von <code>m3.xlarge</code> Instance-Typen. Dadurch wird Ihre Auswahl an Instance-Typen mit <code>m3.xlarge</code> überschrieben, wodurch Ihre monatlichen Kosten ansteigen könnten.</p>		
resizeClusterMaxInstances	Ein Limit für die maximale Anzahl von Instances, die vom Resize-Algorithmus angefordert werden können.	Ganzzahl
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden.</p> <p>Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.</p>	Aufzählung

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": " myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime

Laufzeitfelder	Description	Slot-Typ
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {" ref": " myRunnableObject Id "}
emrStepLog	Amazon EMR-Schrittprotokolle sind nur bei EMR-Aktivitätsversuchen verfügbar.	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String

Laufzeitfelder	Description	Slot-Typ
@ Zeit healthStatusUpdated	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen.	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen.	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „waitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Sphäre eines Objekts bezeichnet seine Position im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die ein Versuchsobjekt ausführen.	String

Weitere Informationen finden Sie unter:

- [ShellCommandActivity](#)
- [EmrActivity](#)

PigActivity

PigActivity bietet native Unterstützung für Pig-Skripte, AWS Data Pipeline ohne dass die Verwendung von `ShellCommandActivity` oder erforderlich ist `EmrActivity`. PigActivity Unterstützt außerdem Daten-Staging. Wenn das Stage-Feld auf „true“ festgelegt wurde, arrangiert AWS Data Pipeline die Eingabedaten ohne zusätzlichen Code des Benutzers als Schema in Pig.

Beispiel

Im folgenden Pipeline-Beispiel wird gezeigt, wie PigActivity verwendet wird. Die Beispiel-Pipeline führt die folgenden Schritte aus:

- MyPigActivity1 lädt Daten aus Amazon S3 und führt ein Pig-Skript aus, das einige Datenspalten auswählt und sie auf Amazon S3 hochlädt.
- MyPigActivity2 lädt die erste Ausgabe, wählt einige Spalten und drei Datenzeilen aus und lädt sie als zweite Ausgabe auf Amazon S3 hoch.
- MyPigActivity3 lädt die zweiten Ausgabedaten, fügt zwei Datenzeilen und nur die Spalte mit dem Namen „Fifth“ in Amazon RDS ein.

- MyPigActivity4 lädt Amazon RDS-Daten, wählt die erste Datenzeile aus und lädt sie auf Amazon S3 hoch.

```
{  
  "objects": [  
    {  
      "id": "MyInputData1",  
      "schedule": {  
        "ref": "MyEmrResourcePeriod"  
      },  
      "directoryPath": "s3://amzn-s3-demo-bucket/pigTestInput",  
      "name": "MyInputData1",  
      "dataFormat": {  
        "ref": "MyInputDataType1"  
      },  
      "type": "S3DataNode"  
    },  
    {  
      "id": "MyPigActivity4",  
      "scheduleType": "CRON",  
      "schedule": {  
        "ref": "MyEmrResourcePeriod"  
      },  
      "input": {  
        "ref": "MyOutputData3"  
      },  
      "pipelineLogUri": "s3://amzn-s3-demo-bucket/path/",  
      "name": "MyPigActivity4",  
      "runsOn": {  
        "ref": "MyEmrResource"  
      },  
      "type": "PigActivity",  
      "dependsOn": {  
        "ref": "MyPigActivity3"  
      },  
      "output": {  
        "ref": "MyOutputData4"  
      },  
      "script": "B = LIMIT ${input1} 1; ${output1} = FOREACH B GENERATE one;",  
      "stage": "true"  
    },  
    {
```

```
"id": "MyPigActivity3",
"scheduleType": "CRON",
"schedule": {
    "ref": "MyEmrResourcePeriod"
},
"input": {
    "ref": "MyOutputData2"
},
"pipelineLogUri": "s3://amzn-s3-demo-bucket/path",
"name": "MyPigActivity3",
"runsOn": {
    "ref": "MyEmrResource"
},
"script": "B = LIMIT ${input1} 2; ${output1} = FOREACH B GENERATE Fifth;",
"type": "PigActivity",
"dependsOn": {
    "ref": "MyPigActivity2"
},
"output": {
    "ref": "MyOutputData3"
},
"stage": "true"
},
{
"id": "MyOutputData2",
"schedule": {
    "ref": "MyEmrResourcePeriod"
},
"name": "MyOutputData2",
"directoryPath": "s3://amzn-s3-demo-bucket/PigActivityOutput2",
"dataFormat": {
    "ref": "MyOutputDataType2"
},
"type": "S3DataNode"
},
{
"id": "MyOutputData1",
"schedule": {
    "ref": "MyEmrResourcePeriod"
},
"name": "MyOutputData1",
"directoryPath": "s3://amzn-s3-demo-bucket/PigActivityOutput1",
"dataFormat": {
    "ref": "MyOutputDataType1"
}
```

```
},
  "type": "S3DataNode"
},
{
  "id": "MyInputDataType1",
  "name": "MyInputDataType1",
  "column": [
    "First STRING",
    "Second STRING",
    "Third STRING",
    "Fourth STRING",
    "Fifth STRING",
    "Sixth STRING",
    "Seventh STRING",
    "Eighth STRING",
    "Ninth STRING",
    "Tenth STRING"
  ],
  "inputRegEx": "^(\\\\\\\$+) (\\\\\\\$+) (\\\\\\\$+) (\\\\\\\$+) (\\\\\\\$+) (\\\\\\\$+) (\\\\\\\$+) (\\\\\$+) (\\\\\\\$+) (\\\\\\\$+)",
  "type": "RegEx"
},
{
  "id": "MyEmrResource",
  "region": "us-east-1",
  "schedule": {
    "ref": "MyEmrResourcePeriod"
  },
  "keyPair": "example-keypair",
  "masterInstanceType": "m1.small",
  "enableDebugging": "true",
  "name": "MyEmrResource",
  "actionOnTaskFailure": "continue",
  "type": "EmrCluster"
},
{
  "id": "MyOutputDataType4",
  "name": "MyOutputDataType4",
  "column": "one STRING",
  "type": "CSV"
},
{
  "id": "MyOutputData4",
  "schedule": {
```

```
        "ref": "MyEmrResourcePeriod"
    },
    "directoryPath": "s3://amzn-s3-demo-bucket/PigActivityOutput3",
    "name": "MyOutputData4",
    "dataFormat": {
        "ref": "MyOutputDataType4"
    },
    "type": "S3DataNode"
},
{
    "id": "MyOutputDataType1",
    "name": "MyOutputDataType1",
    "column": [
        "First STRING",
        "Second STRING",
        "Third STRING",
        "Fourth STRING",
        "Fifth STRING",
        "Sixth STRING",
        "Seventh STRING",
        "Eighth STRING"
    ],
    "columnSeparator": "*",
    "type": "Custom"
},
{
    "id": "MyOutputData3",
    "username": "___",
    "schedule": {
        "ref": "MyEmrResourcePeriod"
    },
    "insertQuery": "insert into #{table} (one) values (?)",
    "name": "MyOutputData3",
    "*password": "___",
    "runsOn": {
        "ref": "MyEmrResource"
    },
    "connectionString": "jdbc:mysql://example-database-instance:3306/example-database",
    "selectQuery": "select * from #{table}",
    "table": "example-table-name",
    "type": " MySqlDataNode"
},
{
```

```
"id": "MyOutputDataType2",
"name": "MyOutputDataType2",
"column": [
    "Third STRING",
    "Fourth STRING",
    "Fifth STRING",
    "Sixth STRING",
    "Seventh STRING",
    "Eighth STRING"
],
"type": "TSV"
},
{
"id": "MyPigActivity2",
"scheduleType": "CRON",
"schedule": {
    "ref": "MyEmrResourcePeriod"
},
"input": {
    "ref": "MyOutputData1"
},
"pipelineLogUri": "s3://amzn-s3-demo-bucket/path",
"name": "MyPigActivity2",
"runsOn": {
    "ref": "MyEmrResource"
},
"dependsOn": {
    "ref": "MyPigActivity1"
},
"type": "PigActivity",
"script": "B = LIMIT ${input1} 3; ${output1} = FOREACH B GENERATE Third, Fourth, Fifth, Sixth, Seventh, Eighth;",
"output": {
    "ref": "MyOutputData2"
},
"stage": "true"
},
{
"id": "MyEmrResourcePeriod",
"startDateTime": "2013-05-20T00:00:00",
"name": "MyEmrResourcePeriod",
"period": "1 day",
"type": "Schedule",
"endDateTime": "2013-05-21T00:00:00"
```

```

},
{
  "id": "MyPigActivity1",
  "scheduleType": "CRON",
  "schedule": {
    "ref": "MyEmrResourcePeriod"
  },
  "input": {
    "ref": "MyInputData1"
  },
  "pipelineLogUri": "s3://amzn-s3-demo-bucket/path",
  "scriptUri": "s3://amzn-s3-demo-bucket/script/pigTestScript.q",
  "name": "MyPigActivity1",
  "runsOn": {
    "ref": "MyEmrResource"
  },
  "scriptVariable": [
    "column1=First",
    "column2=Second",
    "three=3"
  ],
  "type": "PigActivity",
  "output": {
    "ref": "MyOutputData1"
  },
  "stage": "true"
}
]
}

```

Der Inhalt von `pigTestScript.q` ist wie folgt:

```
B = LIMIT ${input1} $three; ${output1} = FOREACH B GENERATE $column1, $column2, Third,
Fourth, Fifth, Sixth, Seventh, Eighth;
```

Syntax

Objektaufruf-Felder	Description	Slot-Typ
schedule	Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Benutzer müssen einen Zeitplanverweis auf ein anderes	Referenzobjekt, zum Beispiel „schedule“:

Objektaufruf-Felder	Description	Slot-Typ
	<p>Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Benutzer können diese Anforderung erfüllen, indem sie explizit einen Zeitplan für das Objekt festlegen, z. B. indem sie „schedule“: {"ref": "DefaultSchedule"} angeben. In den meisten Fällen ist es besser, den Zeitplan zuweisen auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Benutzer ein übergeordnetes Objekt mit Zeitplänenreferenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	{"ref": "myScheduleId"}

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
script	Das auszuführende Pig-Skript.	String
scriptUri	Der Speicherort des auszuführenden Pig-Skripts (z. B. s3://scriptLocation).	String

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
runsOn	EMR-Cluster, auf dem das PigActivity läuft.	Referenzobjekt, zum Beispiel „runsOn“: {"ref": " myEmrClusterId "}
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird ignoriert.workerGroup	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Der zuletzt gemeldete Status von der Remote-Aktivität.	String
attemptTimeout	Das Timeout für die Fertigstellung der Remote-Arbeit. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
dependsOn	Gibt die Abhängigkeit von einem anderen ausführbaren Objekt an.	Referenzobjekt, zum Beispiel „dependSon“: {"ref": " myActivityId „}
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden.	Aufzählung
input	Die Eingangsdatenquelle.	Referenzobjekt, zum Beispiel „input“:

Optionale Felder	Description	Slot-Typ
		{"ref": " myDataNode Id "}
lateAfterTimeout	Die verstrichene Zeit nach dem Start der Pipeline, innerhalb derer das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Die maximale Anzahl von Versuchen bei Ausfällen.	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, zum Beispiel „onFail“: {"ref": " myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, zum Beispiel "onLateAction": {"ref": " myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, zum Beispiel „onSuccess“: {"ref": " myActionId „}
output	Die Eingangsdatenquelle.	Referenzobjekt, zum Beispiel „output“: {"ref": " myDataNode Id "}

Optionale Felder	Description	Slot-Typ
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, zum Beispiel „parent“: {"ref": " myBaseObject Id "}
pipelineLogUri	Die Amazon S3 S3-URI (z. B. 's3://BucketName/Key/ ') für das Hochladen von Protokollen für die Pipeline.	String
postActivityTaskConfig	Post-Activity-Konfigurationsskript, das ausgeführt werden soll. Dies besteht aus einer URI des Shell-Skripts in Amazon S3 und einer Liste von Argumenten.	Referenzobjekt, zum Beispiel "postActivityTaskConfig": {"ref": " myShellScript ConfigId „}
preActivityTaskConfig	Pre-Activity-Konfigurationsskript, das ausgeführt werden soll. Dieses besteht aus einer URI des Shell-Skripts in Amazon S3 und einer Liste von Argumenten.	Referenzobjekt, zum Beispiel "preActivityTaskConfig": {"ref": " myShellScript ConfigId „}
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, zum Beispiel „precondition“: {"ref": " myPreconditionId „}
reportProgressTimeout	Das Timeout für aufeinanderfolgende Aufrufe von <code>reportProgress</code> durch Remote-Arbeit. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum

Optionale Felder	Description	Slot-Typ
resizeClusterBeforeRunning	Ändern Sie die Größe des Clusters, bevor Sie diese Aktivität ausführen, um DynamoDB-Datenknoten aufzunehmen, die als Eingaben oder Ausgaben angegeben sind.	Boolesch
	<p>Note</p> <p>Wenn Ihre Aktivität <code>a</code> entweder <code>DynamoDBDataNode</code> als Eingabe- oder Ausgabedatenknoten verwendet und Sie den Wert auf <code>setzenTRUE</code>, AWS Data Pipeline beginnt die <code>resizeClusterBeforeRunning</code> Verwendung von <code>m3.xlarge</code> Instance-Typen. Dadurch wird Ihre Auswahl an Instance-Typen mit <code>m3.xlarge</code> überschrieben, wodurch Ihre monatlichen Kosten ansteigen könnten.</p>	
resizeClusterMaxInstances	Ein Limit für die maximale Anzahl von Instances, die vom Resize-Algorithmus angefordert werden können.	Ganzzahl
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Zeitreihe nstilplanung bedeutet, dass Instances am Ende jedes Intervalls geplant werden und Cron-Stil-Planung bedeutet, dass Instances zu Beginn jedes Intervalls geplant werden.</p> <p>Ein On-Demand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene scheduleType sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den ActivatePipeline Vorgang einfach für jeden nachfolgenden Lauf auf. Die Werte sind: cron, ondemand und timeseries.</p>	Aufzählung
scriptVariable	Die Argumente, die an das Pig-Skript übergeben werden sollen. Sie können scriptVariable mit script oder scriptUri verwenden.	String
stage	Bestimmt, ob das Staging aktiviert ist, und ermöglicht Ihrem Pig-Skript den Zugriff auf die Tabellen mit Staging-Daten, wie \${INPUT1} und \${OUTPUT1}	Boolesch

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, zum Beispiel „ActiveInstances“: {"ref": "Id"} myRunnableObject
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, zum Beispiel "cascadeFailedOn": {"ref": "myRunnableObject Id"}
emrStepLog	Amazon EMR-Schrittprotokolle sind nur bei EMR-Aktivitätsversuchen verfügbar.	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String

Laufzeitfelder	Description	Slot-Typ
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ Zeit healthStatusUpdated	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für das Objekt einplanen.	DateTime
@scheduledStartTime	Startzeit für das Objekt einplanen.	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde	String

Laufzeitfelder	Description	Slot-Typ
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, zum Beispiel „WaitingOn“: {"ref": "myRunnableObject Id "}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Weitere Informationen finden Sie unter:

- [ShellCommandActivity](#)
- [EmrActivity](#)

RedshiftCopyActivity

Kopiert Daten von DynamoDB oder Amazon S3 nach Amazon Redshift. Sie können Daten in eine neue Tabelle laden oder Daten in einer vorhandenen Tabelle einfach zusammenführen.

Hier finden Sie eine Übersicht über einen Anwendungsfall, in dem RedshiftCopyActivity verwendet wird:

1. Verwenden Sie zunächst AWS Data Pipeline , um Ihre Daten in Amazon S3 bereitzustellen.
2. Wird verwendetRedshiftCopyActivity, um die Daten von Amazon RDS und Amazon EMR nach Amazon Redshift zu verschieben.

Auf diese Weise können Sie Ihre Daten in Amazon Redshift laden, wo Sie sie analysieren können.

3. Wird verwendet [SqlActivity](#), um SQL-Abfragen für die Daten durchzuführen, die Sie in Amazon Redshift geladen haben.

Darüber hinaus unterstützt RedshiftCopyActivity Ihre Arbeit mit einem S3DataNode, weil es eine Manifestdatei unterstützt. Weitere Informationen finden Sie unter [S3 DataNode](#).

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

Um die Formatkonvertierung sicherzustellen, verwendet dieses Beispiel [EMPTYASNULL](#) und [IGNOREBLANKLINES](#), spezielle Konvertierungsparameter in commandOptions. Weitere Informationen finden Sie unter [Datenkonvertierungsparameter](#) im Amazon Redshift Database Developer Guide.

```
{  
  "id" : "S3ToRedshiftCopyActivity",  
  "type" : "RedshiftCopyActivity",  
  "input" : { "ref": "MyS3DataNode" },  
  "output" : { "ref": "MyRedshiftDataNode" },  
  "insertMode" : "KEEP_EXISTING",  
  "schedule" : { "ref": "Hour" },  
  "runsOn" : { "ref": "MyEc2Resource" },  
  "commandOptions": ["EMPTYASNULL", "IGNOREBLANKLINES"]  
}
```

Die folgende Pipeline-Beispieldefinition zeigt eine Aktivität, die den Einfügemodus APPEND nutzt:

```
{  
  "objects": [  
    {  
      "id": "CSVId1",  
      "name": "DefaultCSV1",  

```

```
"name": "DefaultRedshiftDatabase1",
 "*password": "password",
 "type": "RedshiftDatabase",
 "clusterId": "redshiftclusterId"
},
{
  "id": "Default",
  "scheduleType": "timeseries",
  "failureAndRerunMode": "CASCADE",
  "name": "Default",
  "role": "DataPipelineDefaultRole",
  "resourceRole": "DataPipelineDefaultResourceRole"
},
{
  "id": "RedshiftDataNodeId1",
  "schedule": {
    "ref": "ScheduleId1"
  },
  "tableName": "orders",
  "name": "DefaultRedshiftDataNode1",
  "createTableSql": "create table StructuredLogs (requestBeginTime CHAR(30)
PRIMARY KEY DISTKEY SORTKEY, requestEndTime CHAR(30), hostname CHAR(100), requestData
varchar(20));",
  "type": "RedshiftDataNode",
  "database": {
    "ref": "RedshiftDatabaseId1"
  }
},
{
  "id": "Ec2ResourceId1",
  "schedule": {
    "ref": "ScheduleId1"
  },
  "securityGroups": "MySecurityGroup",
  "name": "DefaultEc2Resource1",
  "role": "DataPipelineDefaultRole",
  "logUri": "s3://myLogs",
  "resourceRole": "DataPipelineDefaultResourceRole",
  "type": "Ec2Resource"
},
{
  "id": "ScheduleId1",
  "startDateTime": "yyyy-mm-ddT00:00:00",
  "name": "DefaultSchedule1",
```

```

    "type": "Schedule",
    "period": "period",
    "endDateTime": "yyyy-mm-ddT00:00:00"
},
{
  "id": "S3DataNodeId1",
  "schedule": {
    "ref": "ScheduleId1"
  },
  "filePath": "s3://datapipeline-us-east-1/samples/hive-ads-samples.csv",
  "name": "DefaultS3DataNode1",
  "dataFormat": {
    "ref": "CSVId1"
  },
  "type": "S3DataNode"
},
{
  "id": "RedshiftCopyActivityId1",
  "input": {
    "ref": "S3DataNodeId1"
  },
  "schedule": {
    "ref": "ScheduleId1"
  },
  "insertMode": "APPEND",
  "name": "DefaultRedshiftCopyActivity1",
  "runsOn": {
    "ref": "Ec2ResourceId1"
  },
  "type": "RedshiftCopyActivity",
  "output": {
    "ref": "RedshiftDataNodeId1"
  }
}
]
}

```

Der APPEND-Vorgang fügt Elemente zu einer Tabelle hinzu, unabhängig von Primär- oder Sortierschlüsseln. Bei der folgenden Tabelle können Sie beispielsweise einen Datensatz mit demselben ID- und Benutzer-Wert anfügen.

ID(PK)	USER
1	aaa

2	bbb
---	-----

Sie können einen Datensatz mit demselben ID- und Benutzer-Wert anfügen:

ID(PK)	USER
1	aaa
2	bbb
1	aaa

Note

Wenn ein APPEND-Vorgang unterbrochen und wieder aufgenommen wird, ist es möglich, dass die entstandene Wiederausführungs-Pipeline von Anfang an Anfügungen vornimmt. Dies kann zu weiteren Duplikierungen führen. Sie sollten dieses Verhalten kennen, besonders, wenn Sie Logik verwenden, die die Anzahl an Zeilen zählt.

Ein Tutorial finden Sie unter [Daten mithilfe von Amazon Redshift nach Amazon Redshift kopieren AWS Data Pipeline](#).

Syntax

Pflichtfelder	Description	Slot-Typ
insertMode	<p>Legt fest, AWS Data Pipeline was mit bereits vorhandenen Daten in der Zieltabelle geschehen soll, die sich mit Zeilen in den zu ladenden Daten überschneiden.</p> <p>Gültige Werte sind: KEEP_EXISTING , OVERWRITE_EXISTING , TRUNCATE und APPEND.</p> <p>KEEP_EXISTING fügt der Tabelle neue Zeilen hinzu und lässt die vorhandenen Zeilen unverändert.</p> <p>KEEP_EXISTING und OVERWRITE_EXISTING verwenden den Primärsch</p>	Aufzählung

Pflichtfelder	Description	Slot-Typ
	<p>lüssel, Sortier- und Verteilschlüssel, um zu identifizieren, welche eingehende Zeilen mit vorhandenen Zeilen übereinstimmen. Weitere Informationen finden Sie unter Aktualisieren und Einfügen neuer Daten im Amazon Redshift Database Developer Guide.</p> <p>TRUNCATE löscht alle Daten in der Zieltabelle, bevor die neuen Daten hinzugefügt werden.</p> <p>APPEND fügt alle Datensätze am Ende der Redshift-Tabelle an. APPEND setzt keinen Primär-, Verteilungs- oder Sortierschlüssel voraus. Es können also Zeilen hinzugefügt werden, bei denen es sich um potenzielle Duplikate handelt.</p>	

Objektaufruf-Felder	Description	Slot-Typ
schedule	<p>Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen.</p> <p>Sie müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen.</p> <p>In den meisten Fällen empfehlen wir, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Sie können beispielsweise einen Zeitplan explizit für das Objekt festlegen, indem Sie "schedule": {"ref": "DefaultSchedule"} angeben.</p>	Referenzobjekt, wie z. B.: "schedule": {"ref": "myScheduleId"}

Objektaufruf-Felder	Description	Slot-Typ
	<p>Wenn der Hauptplan in Ihrer Pipeline verschachtelte Zeitpläne enthält, erstellen Sie ein übergeordnetes Objekt mit Zeitplanreferenz.</p> <p>Weitere Informationen zu optionalen Zeitplankonfigurationen finden Sie unter Zeitplan.</p>	

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
runsOn	Die Rechenressource zum Ausführen der Aktivität oder des Befehls. Zum Beispiel eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId"}
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird workerGroup ignoriert.	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum

Optionale Felder	Description	Slot-Typ
commandOptions	<p>Verwendet Parameter, die während des COPY-Vorgangs an den Amazon Redshift Redshift-Datenknoten übergeben werden. Informationen zu Parametern finden Sie unter COPY im Amazon Redshift Database Developer Guide.</p> <p>Wenn COPY die Tabelle lädt, versucht der Befehl implizit, die Zeichenfolgen in den Quelldaten in den Datentyp der Zielspalte zu konvertieren. Zusätzlich zu den Standard-Datenkonvertierungen, die automatisch stattfinden, wenn Fehler erhalten oder andere Konvertierungen benötigen, können Sie zusätzliche Umrechnungsparameter angeben. Weitere Informationen finden Sie unter Datenkonvertierungsparameter im Amazon Redshift Database Developer Guide.</p> <p>Wenn dem Eingabe- oder Ausgabedatenknoten ein Datenformat zugeordnet ist, werden die angegebenen Parameter ignoriert.</p> <p>Da beim Kopieren die Daten zunächst mit dem Befehl COPY in eine Staging-Tabelle eingefügt und danach mit dem Befehl INSERT von der Staging- in die Zieltabelle kopiert werden, können einige COPY-Parameter nicht verwendet werden (z. B. die Fähigkeit des COPY-Befehls, der das automatische Komprimieren der Tabelle aktiviert). Wenn die Tabelle komprimiert werden soll, fügen Sie der Anweisung CREATE TABLE Angaben zur Spaltencodierung hinzu.</p> <p>In einigen Fällen, in denen Daten aus dem Amazon Redshift-Cluster entladen und Dateien in Amazon S3 erstellt werden müssen, ist</p>	String

Optionale Felder	Description	Slot-Typ
	<p>das außerdem auf den UNLOAD Betrieb von Amazon Redshift RedshiftCopyActivity angewiesen.</p> <p>Zur Verbesserung der Leistung beim Kopieren und Entladen geben Sie den PARALLEL OFF-Parameter aus dem UNLOAD Befehl an. Informationen zu Parametern finden Sie unter UNLOAD im Amazon Redshift Database Developer Guide.</p>	
dependsOn	Angeben der Abhängigkeit von einem anderen ausführbaren Objekt.	Referenzobjekt: "dependsOn": {"ref": "myActivityId"}
failureAndRerunMode	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden	Aufzählung
input	Der Eingabedatenknoten. Die Datenquelle kann Amazon S3, DynamoDB oder Amazon Redshift sein.	Referenzobjekt: "input": {"ref": "myDataNodeId"}
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl

Optionale Felder	Description	Slot-Typ
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt: "onFail": {"ref": "myActionId"}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt: "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt: "onSuccess": {"ref": "myActionId"}
output	Der Ausgabedatenknoten. Der Ausgabesp eicherort kann Amazon S3 oder Amazon Redshift sein.	Referenzobjekt: "output": {"ref": "myDataNodeId"}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt: "parent": {"ref": "myBaseObjectID"}
pipelineLogUri	Die S3-URI (z. B. 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt: "precondition": {"ref": "myPreconditionId"}

Optionale Felder	Description	Slot-Typ
Warteschlange	<p>Entspricht der <code>query_group</code> Einstellung in Amazon Redshift, mit der Sie gleichzeitige Aktivitäten anhand ihrer Platzierung in Warteschlangen zuweisen und priorisieren können.</p> <p>In Amazon Redshift sind bis zu 15 gleichzeitige Verbindungen möglich. Weitere Informationen finden Sie unter Zuweisen von Abfragen zu Warteschlangen im Amazon RDS Database Developer Guide.</p>	String
reportProgressTimeout	<p>Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in <code>reportProgress</code>.</p> <p>Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.</p>	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dieser Option können Sie angeben, ob der Plan für die Objekte in Ihrer Pipeline vorgesehen ist. Werte sind <code>cron</code>, <code>ondemand</code> und <code>timeseries</code>.</p> <p>Die <code>timeseries</code> Planung bedeutet, dass Instances am Ende jedes Intervalls geplant sind.</p> <p>Die <code>Cron</code> Planung bedeutet, dass Instances am Anfang jedes Intervalls geplant sind.</p> <p>Ein <code>ondemand</code>-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen.</p> <p>Um <code>ondemand</code>-Pipelines zu verwenden, rufen Sie einfach den <code>ActivatePipeline</code> - Vorgang für jeden nachfolgenden Lauf auf.</p> <p>Wenn Sie einen <code>ondemand</code>-Zeitplan verwenden, müssen Sie ihn im Standardobjekt angeben, und er muss der einzige für die Objekte in der Pipeline angegebene <code>scheduleType</code> sein.</p>	Aufzählung

Optionale Felder	Description	Slot-Typ
transformSql	<p>Der zum Transformieren der Eingabedaten verwendete SQL SELECT-Ausdruck.</p> <p>Führen Sie den Ausdruck <code>transformSql</code> in der Tabelle mit dem Namen <code>staging</code> aus.</p> <p>Wenn Sie Daten aus DynamoDB oder Amazon S3 kopieren, AWS Data Pipeline erstellt eine Tabelle namens „Staging“ und lädt zunächst Daten hinein. Die Daten dieser Tabelle werden zum Aktualisieren der Zieltabelle verwendet.</p> <p>Das Ausgabe-Schema von <code>transformSql</code> muss mit dem Schema der endgültigen Zieltabelle übereinstimmen.</p> <p>Wenn Sie die Option <code>transformSql</code> angeben, wird von der angegebenen SQL-Anweisung eine zweite Staging-Tabelle erstellt. Die Daten dieser zweiten Staging-Tabelle werden anschließend in die endgültige Zieltabelle übernommen.</p>	String

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt: "activeInstances": {"ref": "myRunnableObjectInstanceId"}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime

Laufzeitfelder	Description	Slot-Typ
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt: "cascadeFailedOn": {"ref": "myRunnable ObjectId"}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String

Laufzeitfelder	Description	Slot-Typ
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ Zeit healthStatusUpdated	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen.	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen.	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt: "waitingOn": {"ref":"myRunnableObject.Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Sphäre eines Objekts. Gibt seine Position im Lebenszyklus an. Beispielsweise ergeben Komponentenobjekte Instance-Objekte, die Versuchsobjekte ausführen.	String

ShellCommandActivity

Führt einen Befehl oder ein Skript aus. Mit `ShellCommandActivity` können Sie Zeitreihen oder Cron-ähnliche geplante Aufgaben ausführen.

Wenn das `stage` Feld auf `true` gesetzt ist und mit einem verbunden wird `S3DataNode`, `ShellCommandActivity` unterstützt es das Konzept der Datenbereitstellung, was bedeutet, dass Sie Daten von Amazon S3 an einen Staging-Speicherort, wie Amazon EC2 oder Ihre lokale Umgebung, verschieben können, die Daten mithilfe von Skripten und dem `ShellCommandActivity` bearbeiten und sie zurück zu Amazon S3 verschieben können.

Wenn in diesem Fall Ihr Shell-Befehl mit einem Eingabe-`S3DataNode` verbunden ist, werden Ihre Shell-Skripts mit `${INPUT1_STAGING_DIR}`, `${INPUT2_STAGING_DIR}` und anderen Feldern ausgeführt, die auf die `ShellCommandActivity`-Eingabefelder verweisen.

In ähnlicher Weise kann die Ausgabe des Shell-Befehls in einem Ausgabeverzeichnis bereitgestellt werden, um automatisch an Amazon S3 weitergeleitet zu werden, auf das mit, verwiesen wird `${OUTPUT1_STAGING_DIR} ${OUTPUT2_STAGING_DIR}`, usw.

Diese Ausdrücke können als Befehlszeilenargumente zum Shell-Befehl weitergeleitet werden, sodass Sie sie für Datentransformationslogik verwenden können.

`ShellCommandActivity` gibt Linux-ähnliche Fehlercodes und Zeichenfolgen aus. Wenn `ShellCommandActivity` fehlschlägt, ist der angezeigte `error` ein Wert ungleich Null.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
  "id" : "CreateDirectory",
  "type" : "ShellCommandActivity",
  "command" : "mkdir new-directory"
}
```

Syntax

Objektaufruf-Felder	Description	Slot-Typ
schedule	<p>Dieses Objekt wird innerhalb der Ausführung eines schedule-Intervalls aufgerufen.</p> <p>Um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen, geben Sie eine schedule-Referenz auf ein anderes Objekt an.</p> <p>Um diese Anforderung zu erfüllen, setzen Sie explizit einen schedule auf das Objekt, z. B. mit "schedule": {"ref": "DefaultSchedule"} .</p> <p>In den meisten Fällen ist es besser, die schedule-Referenz auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Besteht die Pipeline aus einem Baum mit Zeitplänen (Zeitpläne innerhalb des Hauptplans), erstellen Sie ein übergeordnetes Objekt, das eine Zeitplanreferenz besitzt.</p> <p>Um die Last zu verteilen, AWS Data Pipeline erstellt physische Objekte etwas früher als geplant, führt sie aber termingerecht aus.</p>	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId"}

Objektaufruf-Felder	Description	Slot-Typ
Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html .		
Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)		
command	Den auszuführenden Befehl. Verwenden Sie \$, um auf Positionsparameter zu verweisen , und geben Sie mit scriptArgument die Parameter für den Befehl an. Dieser Wert und alle zugehörigen Parameter müssen in der Umgebung funktionieren, in der Sie den Task-Runner ausführen.	String
scriptUri	Ein Amazon S3-URI-Pfad für eine Datei, die heruntergeladen und als Shell-Befehl ausgeführt werden soll. Geben Sie nur ein Feld scriptUri oder command an. scriptUri kann keine Parameter verwenden. Verwenden Sie stattdessen command.	String

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
runsOn	Die Rechenressource zur Ausführung der Aktivität oder des Befehls, z. B. eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId"}
workerGroup	Wird für Routing-Aufgaben verwendet. Wenn Sie einen <code>runsOn</code> -Wert angeben und <code>workerGroup</code> vorhanden ist, wird ignoriert.	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Der zuletzt gemeldete Status von der Remote-Aktivität.	String
attemptTimeout	Das Timeout für die Fertigstellung der Remote-Arbeit. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
dependsOn	Gibt eine Abhängigkeit von einem anderen ausführbaren Objekt an.	Referenzobjekt, z. B. „dependsOn“: {"ref": "myActivityId"}
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden.	Aufzählung

Optionale Felder	Description	Slot-Typ
input	Der Speicherort der Eingabedaten.	Referenzobjekt, z. B. „input“: {"ref": "myDataNode Id "}
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Die maximale Anzahl von Versuchen bei Ausfällen.	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId „}
output	Der Speicherort der Ausgabedaten.	Referenzobjekt, z. B. „output“: {"ref": "myDataNode Id "}
übergeordneter	Das übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}

Optionale Felder	Description	Slot-Typ
pipelineLogUri	Die Amazon S3 S3-URI, z. B. 's3://BucketName/Key/' für das Hochladen von Protokollen für die Pipeline.	String
precondition	Definiert optional eine Vorbedingung. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. „precondition“: {"ref": "myPreconditionId"}
reportProgressTimeout	Das Timeout für aufeinanderfolgende Aufrufe von <code>reportProgress</code> durch Remote-Aktivitäten. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Gestattet Ihnen, anzugeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen.</p> <p>Die Werte sind: <code>cron</code>, <code>ondemand</code> und <code>timeseries</code>.</p> <p><code>timeseries</code> bedeutet, dass Instances am Ende jedes Intervalls geplant sind.</p> <p><code>Cron</code> bedeutet, dass Instances am Anfang jedes Intervalls geplant sind.</p> <p><code>ondemand</code> bedeutet, Sie können eine Pipeline jeweils einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen <code>ondemand</code>-Zeitplan verwenden, geben Sie ihn im Standardobjekt als einzigen <code>scheduleType</code> für Objekte in der Pipeline an. Um <code>ondemand</code>-Pipelines zu verwenden, rufen Sie einfach den <code>ActivatePipeline</code> -Vorgang für jeden nachfolgenden Lauf auf.</p>	Aufzählung

Optionale Felder	Description	Slot-Typ
scriptArgument	<p>Ein Zeichenfolgenarray im JSON-Format, das dem von dem Befehl angegebenen Befehl übergeben wird. Ist der Befehl beispielsweise echo \$1 \$2, geben Sie scriptArgument als "param1", "param2" an. Für mehrere Argumente und Parameter übergeben Sie das scriptArgument wie folgt: "scriptArgument": "arg1", "scriptArgument": "param1", "scriptArgument": "arg2", "scriptArgument": "param2" Das scriptArgument ist command. Die Verwendung durch einen Fehler scriptUri verursacht.</p>	String
stage	Legt fest, ob Staging aktiviert ist, und gewährt Ihren Shell-Befehlen den Zugriff auf Staging-Datenvariablen, z. B. \${INPUT1_STAGING_DIR} und \${OUTPUT1_STAGING_DIR} .	Boolesch
stderr	Der -Pfad, zu dem Systemfehlermeldungen vom Befehl umgeleitet werden. Wenn Sie das runsOn Feld verwenden, muss es sich um einen Amazon S3 S3-Pfad handeln, da die Ressource, auf der Ihre Aktivität ausgeführt wird, vorübergehend ist. Wenn Sie jedoch das Feld workerGroup angeben, ist ein lokaler Dateipfad zulässig.	String

Optionale Felder	Description	Slot-Typ
stdout	Der Amazon S3 S3-Pfad, der die umgeleitete Ausgabe des Befehls empfängt. Wenn Sie das <code>runOn</code> Feld verwenden, muss es sich um einen Amazon S3 S3-Pfad handeln, da die Ressource, auf der Ihre Aktivität ausgeführt wird, vorübergehend ist. Wenn Sie jedoch das Feld <code>workerGroup</code> angeben, ist ein lokaler Dateipfad zulässig.	String

Laufzeitfelder	Description	Slot-Typ
<code>@activeInstances</code>	Die Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ <code>activeInstances</code> : {"ref": "Id"} myRunnableObject
<code>@actualEndTime</code>	Der Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
<code>@actualStartTime</code>	Der Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
<code>cancellationReason</code>	Der <code>cancellationReason</code> , wenn dieses Objekt storniert wurde.	String
<code>@cascadeFailedOn</code>	Die Beschreibung der Zuständigkeitskette, die den Objektausfall verursacht hat.	Referenzobjekt, z. B. "cascadeFailedOn,: {"ref": "myRunnableObject Id"}
<code>emrStepLog</code>	Amazon EMR-Schrittprotokolle sind nur bei Amazon EMR-Aktivitätsversuchen verfügbar.	String

Laufzeitfelder	Description	Slot-Typ
errorId	Die <code>errorId</code> , wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die <code>errorMessage</code> , wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der das Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle sind bei Versuchen für Amazon EMR-basierte Aktivitäten verfügbar.	String
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Die ID des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ Zeit healthStatusUpdated	Der Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Der Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Der Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime

Laufzeitfelder	Description	Slot-Typ
@nextRunTime	Der Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Die geplante Endzeit für das Objekt.	DateTime
@scheduledStartTime	Die geplante Startzeit für das Objekt.	DateTime
@Status	Der Status des Objekts.	String
@Version	Die AWS Data Pipeline Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Die Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „WaitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Der Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Die Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Position eines Objekts im Lebenszyklus. Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Weitere Informationen finden Sie unter:

- [CopyActivity](#)

- [EmrActivity](#)

SqlActivity

Führt eine SQL-Abfrage (Skript) auf einer Datenbank aus.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
  "id" : "MySqlActivity",
  "type" : "SqlActivity",
  "database" : { "ref": "MyDatabaseID" },
  "script" : "SQLQuery" | "scriptUri" : s3://scriptBucket/query.sql,
  "schedule" : { "ref": "MyScheduleID" },
}
```

Syntax

Pflichtfelder	Description	Slot-Typ
Datenbank	Die Datenbank für die Ausführung des bereitgestellten SQL-Skripts.	Referenzobjekt, z. B. „database“: {"ref": "myDatabaseId"}

Objektaufruf-Felder	Description	Slot-Typ
schedule	Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Sie müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Sie können einen Zeitplan explizit für das Objekt festlegen, indem Sie "schedule": {"ref": "DefaultSchedule"} angeben.	Referenzobjekt, z. B. „schedule“: {"ref": "myScheduleId"}

Objektaufruf-Felder	Description	Slot-Typ
	<p>In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben.</p> <p>Wenn die Pipeline über einen Baum über in den Hauptplan verschachtelte Zeitplänen, können Benutzer ein übergeordnetes Objekt mit Zeitplanreferenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
script	Das auszuführende SQL-Skript. Sie müssen das Skript oder scriptUri angeben. Wenn das Skript in Amazon S3 gespeichert ist, wird das Skript nicht als Ausdruck ausgewertet. Die Angabe mehrerer Werte für scriptArgument ist hilfreich, wenn das Skript in Amazon S3 gespeichert ist.	String
scriptUri	Ein URI, der den Speicherort eines SQL-Skripts angibt, das in dieser Aktivität ausgeführt wird.	String

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
runsOn	Die Rechenressource zum Ausführen der Aktivität oder des Befehls. Zum Beispiel eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster.	Referenzobjekt, z. B. „runsOn“: {"ref": "myResourceId" „}
workerGroup	Die Auftragnehmergruppe. Dies wird für Routing-Aufgaben verwendet. Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird ignoriert.workerGroup	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
dependsOn	Angeben der Abhängigkeit von einem anderen ausführbaren Objekt.	Referenzobjekt, z. B. „dependSon“: {"ref": "myActivityId" „}
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden	Aufzählung

Optionale Felder	Description	Slot-Typ
input	Speicherort der Eingabedaten.	Referenzobjekt, z. B. „input“: {"ref": "myDataNode Id "}
lateAfterTimeout	Der Zeitraum seit dem geplanten Start der Pipeline, in dem die Objektausführung starten muss.	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt in dem durch 'lateAfterTimeout' angegebenen Zeitraum seit dem geplanten Start der Pipeline noch nicht geplant oder immer noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId „}
output	Speicherort der Ausgabedaten. Dies ist nur nützlich, um innerhalb eines Skripts zu referenzieren (z. B. # {output.tablename}) und um die Ausgabetafel zu erstellen, indem 'createTableSql' im Ausgabedatenknoten gesetzt wird. Die Ausgabe der SQL-Abfrage wird nicht in den Ausgabedatenknoten geschrieben.	Referenzobjekt, z. B. „output“: {"ref": "myDataNode Id "}

Optionale Felder	Description	Slot-Typ
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}
pipelineLogUri	Die S3-URI (wie 's3://BucketName/Key/ ') zum Hochladen von Protokollen für die Pipeline.	String
precondition	Legen Sie optional eine Vorbedingung fest. Ein Datenknoten ist solange nicht als "BEREIT" markiert, bis alle Vorbedingungen erfüllt sind.	Referenzobjekt, z. B. „precondition“: {"ref": "myPreconditionId"}
Warteschlange	[Nur Amazon Redshift] Entspricht der Einstellung query_group in Amazon Redshift, mit der Sie gleichzeitig auszuführende Aktivitäten anhand ihrer Platzierung in Warteschlangen zuweisen und priorisieren können. In Amazon Redshift sind bis zu 15 gleichzeitige Verbindungen möglich. Weitere Informationen finden Sie unter Zuweisen von Abfragen zu Warteschlangen im Amazon Redshift Datenbankentwicklungshandbuch.	String
reportProgressTime out	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Werte sind <code>cron</code>, <code>ondemand</code> und <code>timeseries</code>.</p> <p><code>timeseries</code> Planung bedeutet, dass Instances am Ende jedes Intervalls geplant sind.</p> <p><code>cron</code> Planung bedeutet, dass Instances am Anfang jedes Intervalls geplant sind.</p> <p>Ein ondemand-Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Dies bedeutet, dass Sie die Pipeline nicht klonen oder neu erstellen müssen, um sie erneut auszuführen. Wenn Sie einen ondemand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene <code>scheduleType</code> sein. Um ondemand-Pipelines zu verwenden, rufen Sie einfach den <code>ActivatePipeline</code>-Vorgang für jeden nachfolgenden Lauf auf.</p>	Aufzählung
scriptArgument	<p>Eine Liste der Variablen für das Skript. Sie können alternativ Ausdrücke direkt in das Skriptfeld einfügen. Mehrere Werte für <code>scriptArgument</code> sind hilfreich, wenn das Skript in Amazon S3 gespeichert ist. Beispiel:</p> <pre># {format (@scheduledStartTime, "YY-MM-DD HH:MM:SS")}\n# {format (plusPeriod (@scheduledStartTime, „1 Tag“), "HH:MM:SS")} YY-MM-DD</pre>	String

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ActiveInstances“: {"ref": "Id"} myRunnableObject
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {"ref": "myRunnableObject Id"}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String

Laufzeitfelder	Description	Slot-Typ
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ healthStatusUpdated Zeit	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Laufzeitfelder	Description	Slot-Typ
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „waitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Ressourcen

Im Folgenden sind die AWS Data Pipeline Ressourcenobjekte aufgeführt:

Objekte

- [Ec2Resource](#)
- [EmrCluster](#)
- [HttpProxy](#)

Ec2Resource

Eine EC2 Amazon-Instance, die die durch eine Pipeline-Aktivität definierte Arbeit ausführt.

AWS Data Pipeline unterstützt IMDSv2 jetzt die EC2 Amazon-Instance, die eine sitzungsorientierte Methode verwendet, um die Authentifizierung beim Abrufen von Metadateninformationen von

Instances besser handhaben zu können. Eine Sitzung beginnt und beendet eine Reihe von Anfragen, die Software, die auf einer EC2 Amazon-Instance ausgeführt wird, verwendet, um auf die lokal gespeicherten EC2 Amazon-Instance-Metadaten und Anmeldeinformationen zuzugreifen. Die Software startet eine Sitzung mit einer einfachen HTTP-PUT-Anfrage an. IMDSv2 gibt ein geheimes Token an die auf der EC2 Amazon-Instance ausgeführte Software zurück, die das Token als Passwort IMDSv2 für Anfragen nach Metadaten und Anmeldeinformationen verwendet.

Note

Um es IMDSv2 für Ihre EC2 Amazon-Instance zu verwenden, müssen Sie die Einstellungen ändern, da das Standard-AMI nicht kompatibel mit ist IMDSv2. Sie können eine neue AMI-Version angeben, die Sie über den folgenden SSM-Parameter abrufen können:`/aws/service/ami-amazon-linux-latest/amzn-ami-hvm-x86_64-ebs`.

Informationen zu EC2 Standard-Amazon-Instances, die AWS Data Pipeline erstellt werden, wenn Sie keine Instance angeben, finden Sie unter [EC2 Standard-Amazon-Instances nach AWS-Region](#).

Beispiele

EC2-Klassiker

Important

Nur AWS Konten, die vor dem 4. Dezember 2013 erstellt wurden, unterstützen die EC2 - Classic-Plattform. Wenn Sie über eines dieser Konten verfügen, haben Sie möglicherweise die Möglichkeit, EC2 Ressourcenobjekte für eine Pipeline in einem EC2 -Classic-Netzwerk und nicht in einer VPC zu erstellen. Es wird dringend empfohlen, Ressourcen für alle Ihre Pipelines in zu erstellen. VPCs Wenn Sie bereits Ressourcen in EC2 -Classic haben, empfehlen wir Ihnen außerdem, diese auf eine VPC zu migrieren.

Das folgende Beispielobjekt startet eine EC2 Instanz in EC2 -Classic, wobei einige optionale Felder gesetzt sind.

```
{  
  "id" : "MyEC2Resource",  
  "type" : "Ec2Resource",
```

```

"actionOnTaskFailure" : "terminate",
"actionOnResourceFailure" : "retryAll",
"maximumRetries" : "1",
"instanceType" : "m5.large",
"securityGroups" : [
    "test-group",
    "default"
],
"keyPair" : "my-key-pair"
}

```

EC2-VPC

Das folgende Beispielobjekt startet eine EC2 Instanz in einer nicht standardmäßigen VPC, wobei einige optionale Felder gesetzt sind.

```

{
    "id" : "MyEC2Resource",
    "type" : "Ec2Resource",
    "actionOnTaskFailure" : "terminate",
    "actionOnResourceFailure" : "retryAll",
    "maximumRetries" : "1",
    "instanceType" : "m5.large",
    "securityGroupIds" : [
        "sg-12345678",
        "sg-12345678"
    ],
    "subnetId": "subnet-12345678",
    "associatePublicIpAddress": "true",
    "keyPair" : "my-key-pair"
}

```

Syntax

Pflichtfelder	Description	Slot-Typ
resourceRole	Die IAM-Rolle, die die Ressourcen steuert, auf die die EC2 Amazon-Instance zugreifen kann.	String
role	Die IAM-Rolle, mit der die AWS Data Pipeline Instance erstellt wird EC2 .	String

Objektaufruf-Felder	Description	Slot-Typ
schedule	<p>Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen.</p> <p>Sie müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihenfolge für dieses Objekt festzulegen. Sie können dafür eine der folgenden Möglichkeiten auswählen:</p> <ul style="list-style-type: none"> Um sicherzustellen, dass alle Objekte in der Pipeline den Zeitplan übernehmen, legen Sie einen Zeitplan explizit für das Objekt fest: <code>"schedule": {"ref": "DefaultSchedule"}</code>. In den meisten Fällen ist es nützlich, den Zeitplanverweis auf das Standard-Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan übernehmen. Wenn der Hauptplan in Ihrer Pipeline verschachtelte Zeitpläne enthält, können Sie ein übergeordnetes Objekt mit Zeitplanreferenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html. 	Referenzobjekt, z. B. "schedule ": {"ref": "mySchedu leId"}

Optionale Felder	Description	Slot-Typ
actionOnResourceFailure	Die Aktion, die nach einem Ressourcenfehler dieser Ressource ausgeführt wird. Gültige Werte sind "retryall" und "retrynone".	String

Optionale Felder	Description	Slot-Typ
actionOnTaskFehlsc hlag	Die Aktion, die nach einem Aufgabenfehler dieser Ressource ausgeführt wird. Gültige Werte sind "continue" oder "terminat e" .	String
associatePublicIpA dresse	Gibt an, ob der Instance eine öffentliche IP-Adresse zugewiesen wird. Wenn sich die Instance in Amazon EC2 oder Amazon VPC befindet, ist true der Standardwert. Andernfal ls ist der Standardwert false.	Boolesch
attemptStatus	Zuletzt gemeldeter Status von der Remote-Ak tivität.	String
attemptTimeout	Timeout für die Fertigstellung der Remote- Arbeit. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
availabilityZone	Die Availability Zone, in der die EC2 Amazon-In stance gestartet werden soll.	String
deaktivieren IMDSv1	Der Standardwert ist falsch und aktiviert sowohl als auch IMDSv1 IMDSv2. Wenn Sie ihn auf true setzen, wird er deaktiviert IMDSv1 und bietet nur IMDSv2s	Boolesch
failureAndRerunMod us	Beschreibt das Verhalten des Konsument enknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden.	Aufzählung
httpProxy	Der Proxyhost, den Clients verwenden, um eine Verbindung zu AWS Diensten herzustellen.	Referenzobjekt, z. B. "httpProxy": {"ref":"myHttpP roxyId"}

Optionale Felder	Description	Slot-Typ
imageId	Die ID des für die Instance zu verwendenden AMI. AWS Data Pipeline Verwendet standardmäßig den Virtualisierungstyp HVM AMI. Das IDs verwendete spezifische AMI basiert auf einer Region. Sie können das Standard-AMI überschreiben, indem Sie das von Ihnen gewählte HVM-AMI angeben. Weitere Informationen zu AMI-Typen finden Sie unter Linux AMI Virtualization Types und Finding a Linux AMI im EC2 Amazon-Benutzerhandbuch.	String
initTimeout	Die Zeit, die auf den Start der Ressource gewartet wird.	Zeitraum
instanceCount	Als veraltet gekennzeichnet.	Ganzzahl
instanceType	Der Typ der EC2 Amazon-Instance, die gestartet werden soll.	String
keyPair	Der Name des Schlüsselpaars. Wenn Sie eine EC2 Amazon-Instance starten, ohne ein key pair anzugeben, können Sie sich nicht bei ihr anmelden.	String
lateAfterTimeout	Die verstrichene Zeit nach dem Start der Pipeline, innerhalb derer das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Die maximale Anzahl von Versuchen bei Ausfällen.	Ganzzahl

Optionale Felder	Description	Slot-Typ
minInstanceCount	Als veraltet gekennzeichnet.	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. "onFail": {"ref": "myActionId"}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant wurde oder noch ausgeführt wird.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. "onSuccess": {"ref": "myActionId"}
übergeordneter	Das übergeordnete Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. "parent": {"ref": "myBaseObjectID"}
pipelineLogUri	Die Amazon S3 S3-URI (z. B. 's3://BucketName/Key/') zum Hochladen von Protokollen für die Pipeline.	String
Region	Der Code für die Region, in der die EC2 Amazon-Instance ausgeführt werden soll. Standardmäßig wird die Instance in derselben Region wie die Pipeline ausgeführt. Sie können die Instance in derselben Region als abhängiges Datenset ausführen.	Aufzählung

Optionale Felder	Description	Slot-Typ
reportProgressTimeout	Das Timeout für aufeinanderfolgende Aufrufe von <code>reportProgress</code> durch Remote-Arbeit. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum
runAsUser	Der Benutzer, der das ausführen soll TaskRunner.	String
runsOn	Dieses Feld ist für dieses Objekt nicht zulässig.	Referenzobjekt, z. B. "runsOn": {"ref": "myResourceId"}

Optionale Felder	Description	Slot-Typ
scheduleType	<p>Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang oder am Ende eines Intervalls oder bedarfsabhängig geplant werden sollen.</p> <p>Die Werte sind:</p> <ul style="list-style-type: none"> • <code>timeseries</code>. Instanzen werden am Ende jedes Intervalls geplant. • <code>cron</code>. Instanzen werden zu Beginn jedes Intervalls geplant. • <code>ondemand</code>. Ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Sie müssen die Pipeline nicht klonen oder neu erstellen, um sie erneut auszuführen. Wenn Sie einen On-Demand-Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene <code>scheduleType</code> sein. Um On-Demand-Pipelines zu verwenden, rufen Sie den <code>ActivatePipeline</code>-Vorgang für jeden nachfolgenden Lauf auf. 	Aufzählung
securityGroupIds	Die IDs von einer oder mehreren EC2 Amazon-Sicherheitsgruppen, die für die Instances im Ressourcenpool verwendet werden sollen.	String
securityGroups	Eine oder mehrere EC2 Amazon-Sicherheitsgruppen, die für die Instances im Ressource npool verwendet werden sollen.	String

Optionale Felder	Description	Slot-Typ
spotBidPrice	Die maximale Datenmenge pro Stunde für Ihre Spot-Instance in Dollar, wobei es sich um einen Dezimalwert zwischen 0 und einschließlich 20,00 handelt.	String
subnetId	Die ID des EC2 Amazon-Subnetzes, in dem die Instance gestartet werden soll.	String
terminateAfter	Die Anzahl der Stunden, nach denen die Ressource zu beenden ist.	Zeitraum
useOnDemandOnLastAttempt	Dieses Feld bestimmt, ob beim letzten Versuch, eine Spot-Instance anzufordern, stattdessen eine On-Demand-Instance angefordert wird. Auf diese Weise wird sichergestellt, dass wenn die vorherigen Versuche fehlgeschlagen sind, der letzte Versuch nicht unterbrochen wird.	Boolesch
workerGroup	Dieses Feld ist für dieses Objekt nicht zulässig.	String

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. "activeInstances": {"ref": "myRunnableObjectId"}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime

Laufzeitfelder	Description	Slot-Typ	
cancellationReason	Der cancellationReason , wenn dieses Objekt storniert wurde.	String	
@cascadeFailedOn	Beschreibung der Zuständigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn": {"ref": "myRunnable Object", "id": "123"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0", "lastModified": "2012-10-29T12:00:00Z"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0", "lastModified": "2012-10-29T12:00:00Z", "status": "PENDING"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0", "lastModified": "2012-10-29T12:00:00Z", "status": "PENDING", "health": "OK"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0", "lastModified": "2012-10-29T12:00:00Z", "status": "PENDING", "health": "OK", "error": "No error message provided."}	Referenzobjekt, z. B. "cascadeFailedOn": {"ref": "myRunnable Object", "id": "123"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0", "lastModified": "2012-10-29T12:00:00Z"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0", "lastModified": "2012-10-29T12:00:00Z", "status": "PENDING"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0", "lastModified": "2012-10-29T12:00:00Z", "status": "PENDING", "health": "OK"} oder {"ref": "myRunnable Object", "id": "123", "version": "1.0", "lastModified": "2012-10-29T12:00:00Z", "status": "PENDING", "health": "OK", "error": "No error message provided."}
emrStepLog	Schrittprotokolle sind nur bei Amazon EMR-Aktivitätsversuchen verfügbar.	String	
errorId	Die Fehler-ID, wenn dieses Objekt fehlgeschlagen ist.	String	
errorMessage	Die Fehlermeldung, wenn dieses Objekt fehlgeschlagen ist.	String	
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String	
@failureReason	Der Grund für den Ressourcenfehler.	String	
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime	
hadoopJobLog	Hadoop-Jobprotokolle sind bei Versuchen für Amazon EMR-Aktivitäten verfügbar.	String	
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String	
@healthStatusFromInstanceId	Id des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String	

Laufzeitfelder	Description	Slot-Typ
@ Zeit healthStatusUpdated	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Die geplante Endzeit für das Objekt.	DateTime
@scheduledStartTime	Die geplante Startzeit für das Objekt.	DateTime
@Status	Der Status des Objekts.	String
@Version	Die Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. "waitingOn": {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Position eines Objekts im Lebenszyklus. Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

EmrCluster

Stellt die Konfiguration eines Amazon EMR-Clusters dar. Dieses Objekt wird von [EmrActivity](#) und [HadoopActivity](#) zum Starten eines Clusters verwendet.

Inhalt

- [Schedulers](#)
- [Amazon EMR-Release-Versionen](#)
- [Amazon EMR-Berechtigungen](#)
- [Syntax](#)
- [Beispiele](#)
- [Weitere Informationen finden Sie unter:](#)

Schedulers

Scheduler bieten eine Möglichkeit, die Ressourcenzuweisung und Auftragspriorisierung in einem Hadoop-Cluster festzulegen. Administratoren oder Benutzer können einen Scheduler für verschiedene Klassen von Benutzern und Anwendungen auswählen. Ein Scheduler könnte Warteschlangen nutzen, um Ressourcen für Benutzer und Anwendungen zuzuweisen. Sie richten diese Warteschlangen beim Erstellen des Clusters ein. Anschließend können Sie für bestimmte Arbeits- und Benutzertypen eine höhere Priorität festlegen als für andere. Dieses Vorgehen ermöglicht die effiziente Nutzung von Cluster-Ressourcen, wenn mehrere Benutzer Arbeiten zum Cluster übermitteln. Es gibt drei Arten von Schedulern:

- [FairScheduler](#)— Versucht, Ressourcen gleichmäßig über einen längeren Zeitraum einzuplanen.

- [CapacityScheduler](#)— Verwendet Warteschlangen, um Clusteradministratoren zu ermöglichen, Benutzer Warteschlangen mit unterschiedlicher Priorität und Ressourcenzuweisung zuzuweisen.
- Standard: wird vom Cluster verwendet, was über Ihre Site konfiguriert werden kann.

Amazon EMR-Release-Versionen

Eine Amazon-EMR-Version ist eine Gruppe von Open-Source-Anwendungen aus dem Big-Data-Ökosystem. Jede Version umfasst verschiedene Big-Data-Anwendungen, Komponenten und Funktionen, die Sie bei der Erstellung eines Clusters für die Installation und Konfiguration von Amazon EMR auswählen. Sie geben die Version unter Verwendung der Versionsbezeichnung an. Versionsbezeichnungen haben die Form emr-x.x.x. Beispiel, emr-5.30.0. Amazon EMR-Cluster basieren auf dem Release-Label emr-4.0.0 und verwenden später die `releaseLabel` Eigenschaft, um das Release-Label eines `EmrCluster` Objekts anzugeben. Frühere Versionen verwenden die Eigenschaft `amiVersion`.

Important

Alle Amazon EMR-Cluster, die mit Version 5.22.0 oder höher erstellt wurden, verwenden [Signature Version 4](#), um Anfragen an Amazon S3 zu authentifizieren. Einige frühere Versionen verwenden Signature Version 2. Die Unterstützung für Signature Version 2 wird eingestellt. Weitere Informationen finden Sie unter [Amazon S3 Update — SigV2 Deprecation Period Extended and Modified](#). Wir empfehlen dringend, eine Amazon EMR-Release-Version zu verwenden, die Signature Version 4 unterstützt. Für ältere Versionen, beginnend mit EMR 4.7.x, wurde die neueste Version der Serie aktualisiert, um Signature Version 4 zu unterstützen. Wenn Sie eine ältere EMR-Version verwenden, empfehlen wir, die neueste Version der Serie zu verwenden. Vermeiden Sie außerdem Versionen vor EMR 4.7.0.

Überlegungen und Einschränkungen

Verwenden Sie die neueste Version von Task Runner

Wenn Sie ein selbstverwaltetes `EmrCluster` Objekt mit einem Release-Label verwenden, verwenden Sie den neuesten Task Runner. Weitere Informationen zu Task-Runner finden Sie unter [Arbeiten mit Task Runner](#). Sie können Eigenschaftswerte für alle Amazon EMR-Konfigurationsklassifizierungen konfigurieren. Weitere Informationen finden Sie unter [Configuring Applications](#) im Amazon EMR Release Guide, in und in den [the section called “EmrConfiguration” the section called “Eigenschaft” Objektreferenzen](#).

Support für IMDSv2

Bisher nur AWS Data Pipeline unterstützt IMDSv1. AWS Data Pipeline Unterstützt IMDSv2 jetzt Amazon EMR 5.23.1, 5.27.1 und 5.32 oder höher sowie Amazon EMR 6.2 oder höher. IMDSv2 verwendet eine sitzungsorientierte Methode, um die Authentifizierung beim Abrufen von Metadateninformationen von Instances besser handhaben zu können. Sie sollten Ihre Instanzen so konfigurieren, dass sie IMDSv2 Aufrufe tätigen, indem Sie benutzerverwaltete Ressourcen mit -2.0 erstellen. TaskRunner

Amazon EMR 5.32 oder höher und Amazon EMR 6.x

Die Release-Serien Amazon EMR 5.32 oder höher und 6.x verwenden Hadoop Version 3.x, wodurch grundlegende Änderungen bei der Bewertung des Klassenpfads von Hadoop im Vergleich zu Hadoop-Version 2.x eingeführt wurden. Gängige Bibliotheken wie Joda-Time wurden aus dem Klassenpfad entfernt.

Wenn [EmrActivity](#) oder eine [HadoopActivity](#) Jar-Datei ausführt, die Abhängigkeiten von einer Bibliothek hat, die in Hadoop 3.x entfernt wurde, schlägt der Schritt mit dem Fehler oder fehl. `java.lang.NoClassDefFoundError java.lang.ClassNotFoundException` Dies kann bei Jar-Dateien passieren, die mit den Release-Versionen von Amazon EMR 5.x problemlos ausgeführt wurden.

Um das Problem zu beheben, müssen Sie Abhängigkeiten von Jar-Dateien in den Hadoop-Klassenpfad eines `EmrCluster` Objekts kopieren, bevor Sie das oder das starten. `EmrActivity` `HadoopActivity` Dafür stellen wir ein Bash-Skript zur Verfügung. Das Bash-Skript ist an der folgenden Stelle verfügbar, beispielsweise in der AWS Region, in der Ihr `EmrCluster` Objekt ausgeführt wird. `MyRegion us-west-2`

```
s3://datapipeline-MyRegion/MyRegion/bootstrap-actions/latest/TaskRunner/copy-jars-to-hadoop-classpath.sh
```

Die Art und Weise, wie das Skript ausgeführt wird, hängt davon ab, ob `EmrActivity` es auf einer Ressource `HadoopActivity` ausgeführt wird, die von einer selbst verwalteten Ressource verwaltet wird, AWS Data Pipeline oder ob es auf einer selbst verwalteten Ressource ausgeführt wird.

Wenn Sie eine Ressource verwenden, die von verwaltet wird AWS Data Pipeline, fügen Sie dem `EmrCluster` Objekt eine `bootstrapAction` hinzu. Das `bootstrapAction` gibt das Skript und die Jar-Dateien an, die als Argumente kopiert werden sollen. Sie können bis zu 255 `bootstrapAction` Felder pro `EmrCluster` Objekt hinzufügen, und Sie können ein

bootstrapAction Feld zu einem EmrCluster Objekt hinzufügen, das bereits über Bootstrap-Aktionen verfügt.

Um dieses Skript als Bootstrap-Aktion anzugeben, verwenden Sie die folgende Syntax: Dabei JarFileRegion handelt es sich um die Region, in der die Jar-Datei gespeichert ist, und jeweils MyJarFile um den absoluten Pfad einer Jar-Datei in Amazon S3, die in den Hadoop-Klassenpfad kopiert werden soll. Geben Sie standardmäßig keine Jar-Dateien an, die sich im Hadoop-Klassenpfad befinden.

```
s3://datapipeline-MyRegion/MyRegion/bootstrap-actions/latest/TaskRunner/copy-jars-to-hadoop-classpath.sh, JarFileRegion, MyJarFile1, MyJarFile2[, ...]
```

Das folgende Beispiel spezifiziert eine Bootstrap-Aktion, die zwei Jar-Dateien in Amazon S3 kopiert: my-jar-file.jar und dieemr-dynamodb-tool-4.14.0-jar-with-dependencies.jar. Die im Beispiel verwendete Region ist us-west-2.

```
{  
  "id" : "MyEmrCluster",  
  "type" : "EmrCluster",  
  "keyPair" : "my-key-pair",  
  "masterInstanceType" : "m5.xlarge",  
  "coreInstanceType" : "m5.xlarge",  
  "coreInstanceCount" : "2",  
  "taskInstanceType" : "m5.xlarge",  
  "taskInstanceCount": "2",  
  "bootstrapAction" : ["s3://datapipeline-us-west-2/us-west-2/bootstrap-actions/latest/TaskRunner/copy-jars-to-hadoop-classpath.sh, us-west-2, s3://path/to/my-jar-file.jar, s3://dynamodb-dpl-us-west-2/emr-ddb-storage-handler/4.14.0/emr-dynamodb-tools-4.14.0-jar-with-dependencies.jar"]  
}
```

Sie müssen die Pipeline speichern und aktivieren, damit die Änderung an der neuen bootstrapAction Pipeline wirksam wird.

Wenn Sie eine selbstverwaltete Ressource verwenden, können Sie das Skript auf die Clusterinstanz herunterladen und es über die Befehlszeile mit SSH ausführen. Das Skript erstellt ein Verzeichnis mit dem Namen /etc/hadoop/conf/shellprofile.d und eine Datei mit dem Namen datapipeline-jars.sh dieses Verzeichnisses. Die als Befehlszeilenargumente bereitgestellten JAR-Dateien werden in ein Verzeichnis kopiert, das das Skript mit dem Namen erstellt. /home/

`hadoop/datipeline_jars` Wenn Ihr Cluster anders eingerichtet ist, ändern Sie das Skript nach dem Herunterladen entsprechend.

Die Syntax für die Ausführung des Skripts in der Befehlszeile unterscheidet sich geringfügig von der im vorherigen Beispiel `bootstrapAction` gezeigten Syntax. Verwenden Sie Leerzeichen anstelle von Kommas zwischen Argumenten, wie im folgenden Beispiel gezeigt.

```
./copy-jars-to-hadoop-classpath.sh us-west-2 s3://path/to/my-jar-file.jar s3://dynamodb-dpl-us-west-2/emr-ddb-storage-handler/4.14.0/emr-dynamodb-tools-4.14.0-jar-with-dependencies.jar
```

Amazon EMR-Berechtigungen

Wenn Sie eine benutzerdefinierte IAM-Rolle erstellen, sollten Sie sorgfältig die Mindestberechtigungen berücksichtigen, die Ihr Cluster zur Ausführung seiner Aufgaben benötigt. Stellen Sie sicher, dass Sie Zugriff auf die erforderlichen Ressourcen gewähren, z. B. Dateien in Amazon S3 oder Daten in Amazon RDS, Amazon Redshift oder DynamoDB. Wenn Sie `visibleToAllUsers` auf „False“ festlegen möchten, muss Ihre Rolle über die entsprechenden Berechtigungen verfügen. Beachten Sie, dass `DataPipelineDefaultRole` nicht über diese Berechtigungen verfügt. Sie müssen entweder eine Vereinigung der `DataPipelineDefaultRole` Rollen `DefaultDataPipelineResourceRole` und als `EmrCluster` Objektrolle angeben oder zu diesem Zweck Ihre eigene Rolle erstellen.

Syntax

Objektaufruf-Felder	Description	Slot-Typ
<code>schedule</code>	Dieses Objekt wird innerhalb der Ausführung eines Zeitplanintervalls aufgerufen. Sie müssen einen Zeitplanverweis auf ein anderes Objekt angeben, um die Abhängigkeitsausführungsreihefolge für dieses Objekt festzulegen. Sie können diese Anforderung erfüllen, indem Sie explizit einen Zeitplan für das Objekt festlegen, indem sie beispielsweise <code>"schedule": {"ref": "DefaultSchedule"}</code> angeben. In den meisten Fällen ist es besser, den Zeitplanverweis auf das Standard-	Referenzobjekt, z. B. <code>"schedule": {"ref": "myScheduleId"}</code>

Objektaufruf-Felder	Description	Slot-Typ
	<p>Pipeline-Objekt zu setzen, damit alle Objekte diesen Zeitplan erben. Wenn die Pipeline über einen Baum mit Zeitplänen verfügt (Zeitpläne innerhalb des Hauptplans), können Sie ein übergeordnetes Objekt mit Zeitplänenreferenz erstellen. Weitere Informationen zu optionalen Beispiel-Zeitplankonfigurationen finden Sie unter https://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-object-schedule.html.</p>	

Optionale Felder	Description	Slot-Typ
actionOnResourceFehlschlag	Die Aktion, die nach einem Ressourcenfehler dieser Ressource ausgeführt wird. Gültige Werte sind "retryall", der für die festgelegte Dauer wiederholt versucht, alle Aufgaben des Clusters durchzuführen, und "retrynone".	String
actionOnTaskFehlgeschlag	Die Aktion, die nach einem Aufgabenfehler dieser Ressource ausgeführt wird. Gültige Werte sind "continue", was bedeutet, dass der Cluster nicht beendet wird, und "terminate".	String
additionalMasterSecurityGroupIds	Die Kennung zusätzlicher Master-Sicherheitsgruppen des EMR-Clusters, die der Form XXXX6a sg-01 folgt. Weitere Informationen finden Sie unter Zusätzliche Amazon EMR-Sicherheitsgruppen im Amazon EMR Management Guide.	String
additionalSlaveSecurityGroupIds	Die ID zusätzlicher Slave-Sicherheitsgruppen des EMR-Clusters, die dem Format sg-01XXXX6a entspricht.	String

Optionale Felder	Description	Slot-Typ
amiVersion	Die Amazon Machine Image (AMI) -Version, die Amazon EMR zur Installation der Clusterknoten verwendet. Weitere Informationen finden Sie im Amazon EMR-Managementhandbuch .	String
applications	Anwendungen, die im Cluster mit durch Kommas getrennten Argumenten installiert werden sollen. Hive und Pig sind standardmäßig installiert. Dieser Parameter gilt nur für Amazon EMR Version 4.0 und höher.	String
attemptStatus	Der zuletzt gemeldete Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
availabilityZone	Die Availability Zone, in der der Cluster gestartet werden soll.	String
bootstrapAction	Eine Aktion, die ausgeführt werden soll, wenn der Cluster startet. Sie können durch Kommas getrennte Argumente festlegen. Wenn Sie mehrere Aktionen angeben möchten (maximal 255), fügen Sie die entsprechende Anzahl von bootstrapAction -Feldern hinzu. Standardmäßig wird der Cluster ohne Bootstrap-Aktionen gestartet.	String

Optionale Felder	Description	Slot-Typ
Konfiguration	Konfiguration für den Amazon EMR-Cluster. Dieser Parameter gilt nur für Amazon EMR Version 4.0 und höher.	Referenzobjekt, z. B. "configuration": {"ref": "myEmrConfigurationId"}
coreInstanceBidPreis	Der maximale Spot-Preis, den Sie bereit sind, für EC2 Amazon-Instances zu zahlen. Wenn ein Angebotspreis angegeben ist, verwendet Amazon EMR Spot-Instances für die Instance-Gruppe. Angegeben in USD.	String
coreInstanceCount	Gibt an, wie viele Core-Knoten für den Cluster verwendet werden sollen.	Ganzzahl
coreInstanceType	Der Typ der EC2 Amazon-Instance, die für Kernknoten verwendet werden soll. Siehe Unterstützte EC2 Amazon-Instances für Amazon EMR-Cluster .	String
coreGroupConfiguration	Die Konfiguration für die Amazon EMR-Cluster-Core-Instance-Gruppe. Dieser Parameter gilt nur für Amazon EMR Version 4.0 und höher.	Referenzobjekt, z. B. "configuration": {"ref": "myEmrConfigurationId"}
coreEbsConfiguration	Die Konfiguration für Amazon EBS-Volumes, die an jeden der Kernknoten in der Kerngruppe im Amazon EMR-Cluster angehängt werden. Weitere Informationen finden Sie im EC2 Amazon-Benutzerhandbuch unter Instance-Typen, die die EBS-Optimierung Support .	Referenzobjekt, z. B. "coreEbsConfiguration": {"ref": "myEbsConfiguration"}

Optionale Felder	Description	Slot-Typ
customAmiId	Gilt nur für Amazon EMR-Release-Version 5.7.0 und höher. Gibt die AMI-ID eines benutzerdefinierten AMI an, das verwendet werden soll, wenn Amazon EMR EC2 Amazon-Instances bereitstellt. Sie kann auch anstelle von Bootstrap-Aktionen verwendet werden, um Cluster-Knotenkonfigurationen anzupassen. Weitere Informationen finden Sie unter dem folgenden Thema im Amazon EMR Management Guide. Verwenden eines benutzerdefinierten AMI	String
EbsBlockDeviceConfig	<p>Die Konfiguration eines angeforderten Amazon EBS-Blockgeräts, das der Instanzgruppe zugeordnet ist. Diese umfasst eine feste Anzahl an Volumes, die jeder Instance in der Instance-Gruppe zugeordnet wird. Sie umfasst volumesPerInstance und volumeSpecification, wobei:</p> <ul style="list-style-type: none"> • volumesPerInstance die Anzahl der EBS-Volumes mit einer bestimmten Volume-Konfiguration für alle zugeordneten Instances in der Instance-Gruppe ist. • volumeSpecification sind die Amazon EBS-Volume-Spezifikationen, wie Volume-Typ, IOPS und Größe in Gigabytes (GiB), die für das EBS-Volume angefordert werden, das an eine EC2 Instance im Amazon EMR-Cluster angehängt ist. 	Referenzobjekt, z. B. "EbsBlockDeviceConfig": {"ref": "myEbsBlockDeviceConfig"}

Optionale Felder	Description	Slot-Typ
emrManage dMasterSecurityGro upId	Die ID der Master-Sicherheitsgruppe des Amazon EMR-Clusters, die der Form von sg-01XXXX6a folgt. Weitere Informationen finden Sie unter Configure Security Groups im Amazon EMR Management Guide.	String
emrManage dSlaveSecurityGrou pld	Die ID der Slave-Sicherheitsgruppe des Amazon EMR-Clusters, die dem Formular sg-01XXXX6a folgt.	String
enableDebugging	Aktiviert das Debuggen auf dem Amazon EMR-Cluster.	String
failureAndRerunMod us	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden.	Aufzählung
hadoopSchedulerTyp e	Der Scheduler-Typ des Clusters. Gültige Typen sind: PARALLEL_FAIR_SCHEDULING PARALLEL_CAPACITY_SCHEDULING und DEFAULT_SCHEDULER .	Aufzählung
httpProxy	Der Proxy-Host, der von Clients zum Verbinden mit den AWS-Services verwendet wird.	Referenzobjekt, zum Beispiel „HttpProxy“: {"ref": " myHttpProxy Id "}
initTimeout	Die Zeit, die auf den Start der Ressource gewartet wird.	Zeitraum
keyPair	Das EC2 Amazon-Schlüsselpaar, das für die Anmeldung am Master-Knoten des Amazon EMR-Clusters verwendet werden soll.	String

Optionale Felder	Description	Slot-Typ
lateAfterTimeout	Die verstrichene Zeit nach dem Start der Pipeline, innerhalb derer das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
masterInstanceBidPreis	Der maximale Spot-Preis, den Sie bereit sind, für EC2 Amazon-Instances zu zahlen. Es handelt sich um einen Dezimalwert zwischen 0 und einschließlich 20,00. Angegeben in USD. Wenn Sie diesen Wert festlegen, werden Spot-Instances für den Master-Knoten des Amazon EMR-Clusters aktiviert. Wenn ein Angebotspot Preis angegeben ist, verwendet Amazon EMR Spot-Instances für die Instance-Gruppe.	String
masterInstanceType	Der Typ der EC2 Amazon-Instance, die für den Master-Knoten verwendet werden soll. Siehe Unterstützte EC2 Amazon-Instances für Amazon EMR-Cluster .	String
masterGroupConfiguration	Die Konfiguration für die Amazon EMR-Cluster-Master-Instance-Gruppe. Dieser Parameter gilt nur für Amazon EMR Version 4.0 und höher.	Referenzobjekt, z. B. "configuration": {"ref": "myEmrConfigurationId"}
masterEbsConfiguration	Die Konfiguration für Amazon EBS-Volumes, die an jeden der Master-Knoten in der Master-Gruppe im Amazon EMR-Cluster angehängt werden. Weitere Informationen finden Sie im EC2 Amazon-Benutzerhandbuch unter Instance-Typen, die die EBS-Optimierung Support .	Referenzobjekt, z. B. "masterEbsConfiguration": {"ref": "myEbsConfiguration"}

Optionale Felder	Description	Slot-Typ
maxActiveInstances	Die maximale Anzahl gleichzeitiger aktiver Instances einer Komponente. Wiederholungen zählen nicht zur Anzahl der aktiven Instances.	Ganzzahl
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen.	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. "onFail": {"ref": "myActionId"}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. "onSuccess": {"ref": "myActionId"}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. "parent": {"ref": "myBaseObjectID"}
pipelineLogUri	Die Amazon S3 S3-URI (z. B. 's3://BucketName/Key/ ') für das Hochladen von Protokollen für die Pipeline.	String

Optionale Felder	Description	Slot-Typ
Region	Der Code für die Region, in der der Amazon EMR-Cluster ausgeführt werden soll. Standardmäßig wird der Cluster in derselben Region wie die Pipeline ausgeführt. Sie können den Cluster in derselben Region als abhängiges Datenset ausführen.	Aufzählung
releaseLabel	Versionsbezeichnung für den EMR-Cluster	String
reportProgressTime out	Timeout für aufeinanderfolgende Aufrufe von <code>Remote</code> -Arbeit in <code>reportProgress</code> . Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
resourceRole	Die IAM-Rolle, die zur Erstellung des Amazon EMR-Clusters AWS Data Pipeline verwendet wird. Die Standardrolle ist <code>DataPipelineDefaultRole</code> .	String
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum
role	Die IAM-Rolle wurde an Amazon EMR übergeben, um Knoten zu erstellen EC2.	String
runsOn	Dieses Feld ist für dieses Objekt nicht zulässig.	Referenzobjekt, z. B. " <code>runsOn</code> ": {" <code>ref</code> ": " <code>myResourceId</code> "}
Sicherheitskonfiguration	Die ID der EMR-Sicherheitskonfiguration, die auf den Cluster angewendet wird. Dieser Parameter gilt nur für Amazon EMR Version 4.8.0 und höher.	String

Optionale Felder	Description	Slot-Typ
serviceAccessSecurityGroupId	Die ID für die Sicherheitsgruppe für den Servicezugriff des Amazon EMR-Clusters.	Zeichenfolge. Sie hat das Format <code>sg-01XXXX6a</code> , z. B. <code>sg-1234abcd</code> .
scheduleType	Mit dem Zeitplantyp können Sie angeben, ob die Objekte in Ihrer Pipeline-Definition am Anfang des Intervalls oder am Ende des Intervalls geplant werden sollen. Werte sind <code>cron</code> , <code>ondemand</code> und <code>timeseries</code> . Die <code>timeseries</code> -Planung bedeutet, dass Instances am Ende jedes Intervalls geplant sind. Die <code>cron</code> -Planung bedeutet, dass Instances am Anfang jedes Intervalls geplant sind. Ein <code>ondemand</code> -Zeitplan ermöglicht es Ihnen, eine Pipeline einmal pro Aktivierung auszuführen. Sie müssen die Pipeline nicht klonen oder neu erstellen, um sie erneut auszuführen. Wenn Sie einen <code>ondemand</code> -Zeitplan verwenden, muss er im Standardobjekt angegeben werden und der einzige für die Objekte in der Pipeline angegebene <code>scheduleType</code> sein. Um <code>ondemand</code> -Pipelines zu verwenden, rufen Sie einfach den <code>ActivatePipeline</code> -Vorgang für jeden nachfolgenden Lauf auf.	Aufzählung
subnetId	Die ID des Subnetzes, in dem der Amazon EMR-Cluster gestartet werden soll.	String
supportedProducts	Ein Parameter, der Drittanbietersoftware auf einem Amazon EMR-Cluster installiert, z. B. eine Drittanbieter-Distribution von Hadoop.	String

Optionale Felder	Description	Slot-Typ
taskInstanceBidPreis	Der maximale Spot-Preis, den Sie bereit sind, für EC2 Instances zu zahlen. Geben Sie eine Dezimalzahl von 0 bis 20,00 ein. Angegeben in USD. Wenn ein Angebotspreis angegeben ist, verwendet Amazon EMR Spot-Instances für die Instance-Gruppe.	String
taskInstanceCount	Die Anzahl der Task-Knoten, die für den Amazon EMR-Cluster verwendet werden sollen.	Ganzzahl
taskInstanceType	Der Typ der EC2 Amazon-Instance, die für Task-Knoten verwendet werden soll.	String
taskGroupConfiguration	Die Konfiguration für die Amazon EMR-Cluster-Task-Instance-Gruppe. Dieser Parameter gilt nur für Amazon EMR Version 4.0 und höher.	Referenzobjekt, z. B. "configuration": {"ref": "myEmrConfigurationId"}
taskEbsConfiguration	Die Konfiguration für Amazon EBS-Volumes, die an jeden der Task-Knoten in der Aufgabengruppe im Amazon EMR-Cluster angehängt werden. Weitere Informationen finden Sie im EC2 Amazon-Benutzerhandbuch unter Instance-Typen, die die EBS-Optimierung Support .	Referenzobjekt, z. B. "taskEbsConfiguration": {"ref": "myEbsConfiguration"}
terminateAfter	Die Zeitspanne in Stunden, nach der die Ressource beendet wird.	Ganzzahl

Optionale Felder	Description	Slot-Typ
VolumeSpecification	<p>Die Amazon EBS-Volumenspezifikationen, wie Volumetyp, IOPS und Größe in Gigabytes (GiB), die für das Amazon EBS-Volume angefordert werden, das an eine EC2 Amazon-Instance im Amazon EMR-Cluster angehängt ist. Der Knoten kann ein Core-, Master- oder Aufgabenknoten sein.</p> <p>VolumeSpecification enthält:</p> <ul style="list-style-type: none"> • <code>iops()</code> Ganzzahl. Die Anzahl der I/O Operationen pro Sekunde (IOPS), die das Amazon EBS-Volume unterstützt, z. B. 1000. Weitere Informationen finden Sie unter EBS-I/O Characteristics im EC2 Amazon-Benutzerhandbuch. • <code>sizeinGB()</code> . Ganzzahl. Die Größe des Amazon EBS-Volumes in Gibibyte (GiB), zum Beispiel 500. Informationen zu gültigen Kombinationen von Volumetypen und Festplattengrößen finden Sie unter EBS-Volumetypen im EC2 Amazon-Benutzerhandbuch. • <code>volumetType</code> . Schnur. Der Amazon EBS-Volumetyp, zum Beispiel gp2. Es werden die Volume-Typen standard, gp2, io1, st1, sc1 sowie weitere Typen unterstützt. Weitere Informationen finden Sie unter EBS-Volumetypen im EC2 Amazon-Benutzerhandbuch. 	Referenzobjekt, z. B. "VolumeSpecification": {"ref": "myVolume Specification"}

Optionale Felder	Description	Slot-Typ
useOnDemandOnLastAttempt	Dieses Feld bestimmt, ob beim letzten Versuch, eine Ressource anzufordern, eine On-Demand-Instance statt einer Spot-Instance angefordert wird. Auf diese Weise wird sichergestellt, dass wenn die vorherigen Versuche fehlgeschlagen sind, der letzte Versuch nicht unterbrochen wird.	Boolesch
workerGroup	Dieses Feld ist bei diesem Objekt nicht zulässig.	String

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, zum Beispiel „ActiveInstances“: {"ref": "myRunnableObject Id"}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Zuständigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, zum Beispiel "cascadeFailedOn": {"ref": "myRunnableObject Id"}

Laufzeitfelder	Description	Slot-Typ
emrStepLog	Schrittprotokolle sind nur bei Amazon EMR-Aktivitätsversuchen verfügbar.	String
errorId	Die Fehler-ID, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die Fehlermeldung, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
@failureReason	Der Grund für den Ressourcenfehler.	String
@finishedTime	Der Zeitpunkt, zu der dieses Objekt seine Ausführung beendet hat.	DateTime
hadoopJobLog	Hadoop-Jobprotokolle sind bei Versuchen für Amazon EMR-Aktivitäten verfügbar.	String
@healthStatus	Der Integritätsstatus des Objekts, der Erfolg oder Misserfolg der letzten Objekt-Instance widerspiegelt, die einen beendeten Zustand erreicht hat.	String
@healthStatusFromInstanceId	ID des Objekts der letzten Instance, das einen beendeten Zustand erreicht hat.	String
@ Zeit healthStatusUpdated	Zeitpunkt, zu dem der Servicestatus beim letzten Mal aktualisiert wurde.	DateTime
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
@lastDeactivatedTime	Zeitpunkt, zu dem dieses Objekt zuletzt deaktiviert wurde.	DateTime

Laufzeitfelder	Description	Slot-Typ
@ latestCompletedRun Zeit	Zeitpunkt des letzten Laufs, für den die Ausführung abgeschlossen wurde.	DateTime
@latestRunTime	Zeitpunkt des letzten Laufs, für den die Ausführung geplant war.	DateTime
@nextRunTime	Zeitpunkt des Laufs, der als nächstes geplant werden soll	DateTime
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen.	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen.	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, zum Beispiel „WaitingOn“: {"ref": "myRunnableObject Id "}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Position eines Objekts im Lebenszyklus. Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Beispiele

Es folgen Beispiele für diesen Objekttyp.

Inhalt

- [Starten Sie einen Amazon EMR-Cluster mit HadoopVersion](#)
- [Starten Sie einen Amazon EMR-Cluster mit dem Release-Label emr-4.x oder höher](#)
- [Installieren Sie zusätzliche Software auf Ihrem Amazon EMR-Cluster](#)
- [Deaktivieren der serverseitigen Verschlüsselung auf 3.x-Versionen](#)
- [Deaktivieren der serverseitigen Verschlüsselung auf 4.x-Versionen](#)
- [Konfigurieren Sie Hadoop KMS ACLs und erstellen Sie Verschlüsselungszonen in HDFS](#)
- [Festlegen benutzerdefinierter IAM-Rollen](#)
- [Verwenden Sie die EmrCluster Ressource im AWS SDK for Java](#)
- [Einen Amazon EMR-Cluster in einem privaten Subnetz konfigurieren](#)
- [EBS-Volumes zu Cluster-Knoten hinzufügen](#)

Starten Sie einen Amazon EMR-Cluster mit HadoopVersion

Example

Im folgenden Beispiel wird ein Amazon EMR-Cluster mit AMI-Version 1.0 und Hadoop 0.20 gestartet.

```
{  
    "id" : "MyEmrCluster",  
    "type" : "EmrCluster",  
    "hadoopVersion" : "0.20",  
    "keyPair" : "my-key-pair",  
    "masterInstanceType" : "m3.xlarge",  
    "coreInstanceType" : "m3.xlarge",  
    "coreInstanceCount" : "10",  
    "taskInstanceType" : "m3.xlarge",  
    "taskInstanceCount": "10",  
    "bootstrapAction" : ["s3://Region.elasticmapreduce/bootstrap-actions/configure-  
        hadoop,arg1,arg2,arg3","s3://Region.elasticmapreduce/bootstrap-actions/configure-  
        hadoop/configure-other-stuff,arg1,arg2"]  
}
```

Starten Sie einen Amazon EMR-Cluster mit dem Release-Label emr-4.x oder höher

Example

Im folgenden Beispiel wird ein Amazon EMR-Cluster mit dem neueren `releaseLabel` Feld gestartet:

```
{  
  "id" : "MyEmrCluster",  
  "type" : "EmrCluster",  
  "keyPair" : "my-key-pair",  
  "masterInstanceType" : "m3.xlarge",  
  "coreInstanceType" : "m3.xlarge",  
  "coreInstanceCount" : "10",  
  "taskInstanceType" : "m3.xlarge",  
  "taskInstanceCount": "10",  
  "releaseLabel": "emr-4.1.0",  
  "applications": ["spark", "hive", "pig"],  
  "configuration": {"ref": "myConfiguration"}  
}
```

Installieren Sie zusätzliche Software auf Ihrem Amazon EMR-Cluster

Example

`EmrCluster` stellt das `supportedProducts` Feld bereit, das Drittanbieter-Software auf einem Amazon EMR-Cluster installiert. Damit können Sie beispielsweise eine benutzerdefinierte Distribution von Hadoop wie MapR installieren. Er akzeptiert eine durch Kommas getrennte Liste von Argumenten. Die Drittanbieter-Software kann diese Argumente lesen und darauf reagieren. Das folgende Beispiel zeigt, wie Sie mit dem Feld `supportedProducts` von `EmrCluster` einen benutzerdefinierten Cluster der MapR M3-Edition mit Karmasphere Analytics erstellen und ein `EmrActivity`-Objekt darauf ausführen.

```
{  
  "id": "MyEmrActivity",  
  "type": "EmrActivity",  
  "schedule": {"ref": "ResourcePeriod"},  
  "runsOn": {"ref": "MyEmrCluster"},  

```

```
    hdfs://output32113/, -mapper, s3n://elasticmapreduce/samples/wordcount/
    wordSplitter.py, -reducer, aggregate"
},
{
  "id": "MyEmrCluster",
  "type": "EmrCluster",
  "schedule": {"ref": "ResourcePeriod"},
  "supportedProducts": ["mapr,--edition,m3,--version,1.2,--key1,value1", "karmasphere-
  enterprise-utility"],
  "masterInstanceType": "m3.xlarge",
  "taskInstanceType": "m3.xlarge"
}
```

Deaktivieren der serverseitigen Verschlüsselung auf 3.x-Versionen

Example

Eine EmrCluster Aktivität mit einer Hadoop-Version 2.x, erstellt von, AWS Data Pipeline aktiviert standardmäßig serverseitige Verschlüsselung. Wenn Sie die serverseitige Verschlüsselung deaktivieren möchten, müssen Sie eine Bootstrap-Aktion in der Cluster-Objektdefinition festlegen.

Das folgende Beispiel erstellt eine EmrCluster-Aktivität, bei der die serverseitige Verschlüsselung deaktiviert ist:

```
{
  "id": "NoSSEEmrCluster",
  "type": "EmrCluster",
  "hadoopVersion": "2.x",
  "keyPair": "my-key-pair",
  "masterInstanceType": "m3.xlarge",
  "coreInstanceType": "m3.large",
  "coreInstanceCount": "10",
  "taskInstanceType": "m3.large",
  "taskInstanceCount": "10",
  "bootstrapAction": ["s3://Region.elasticmapreduce/bootstrap-actions/configure-
  hadoop, -e, fs.s3.enableServerSideEncryption=false"]
}
```

Deaktivieren der serverseitigen Verschlüsselung auf 4.x-Versionen

Example

Sie müssen die serverseitige Verschlüsselung mit einem EmrConfiguration-Objekt deaktivieren.

Das folgende Beispiel erstellt eine `EmrCluster`-Aktivität, bei der die serverseitige Verschlüsselung deaktiviert ist:

```
{  
    "name": "ReleaseLabelCluster",  
    "releaseLabel": "emr-4.1.0",  
    "applications": ["spark", "hive", "pig"],  
    "id": "myResourceId",  
    "type": "EmrCluster",  
    "configuration": {  
        "ref": "disableSSE"  
    },  
,  
{  
    "name": "disableSSE",  
    "id": "disableSSE",  
    "type": "EmrConfiguration",  
    "classification": "emrfs-site",  
    "property": [{  
        "ref": "enableServerSideEncryption"  
    }  
]  
},  
{  
    "name": "enableServerSideEncryption",  
    "id": "enableServerSideEncryption",  
    "type": "Property",  
    "key": "fs.s3.enableServerSideEncryption",  
    "value": "false"  
}
```

Konfigurieren Sie Hadoop KMS ACLs und erstellen Sie Verschlüsselungszonen in HDFS

Example

Die folgenden Objekte werden ACLs für Hadoop KMS erstellt und erstellen Verschlüsselungszonen und entsprechende Verschlüsselungsschlüssel in HDFS:

```
{  
    "name": "kmsAcls",  
    "id": "kmsAcls",  
    "type": "EmrConfiguration",  
    "classification": "hadoop-kms-acls",
```

```
"property": [
    {"ref":"kmsBlacklist"},
    {"ref":"kmsAcl"}
],
},
{
    "name": "hdfsEncryptionZone",
    "id": "hdfsEncryptionZone",
    "type": "EmrConfiguration",
    "classification": "hdfs-encryption-zones",
    "property": [
        {"ref":"hdfsPath1"},
        {"ref":"hdfsPath2"}
    ]
},
{
    "name": "kmsBlacklist",
    "id": "kmsBlacklist",
    "type": "Property",
    "key": "hadoop.kms.blacklist.CREATE",
    "value": "foo,myBannedUser"
},
{
    "name": "kmsAcl",
    "id": "kmsAcl",
    "type": "Property",
    "key": "hadoop.kms.acl.ROLLOVER",
    "value": "myAllowedUser"
},
{
    "name": "hdfsPath1",
    "id": "hdfsPath1",
    "type": "Property",
    "key": "/myHDFSPath1",
    "value": "path1_key"
},
{
    "name": "hdfsPath2",
    "id": "hdfsPath2",
    "type": "Property",
    "key": "/myHDFSPath2",
    "value": "path2_key"
}
```

Festlegen benutzerdefinierter IAM-Rollen

Example

Wird standardmäßig `DataPipelineDefaultRole` als Amazon EMR-Servicerolle und `DataPipelineDefaultResourceRole` als EC2 Amazon-Instanzprofil AWS Data Pipeline übergeben, um Ressourcen in Ihrem Namen zu erstellen. Sie können jedoch eine benutzerdefinierte Amazon EMR-Servicerolle und ein benutzerdefiniertes Instance-Profil erstellen und diese stattdessen verwenden. AWS Data Pipeline sollte über ausreichende Berechtigungen verfügen, um Cluster mithilfe der benutzerdefinierten Rolle zu erstellen, und Sie müssen sie AWS Data Pipeline als vertrauenswürdige Entität hinzufügen.

Das folgende Beispielobjekt spezifiziert benutzerdefinierte Rollen für den Amazon EMR-Cluster:

```
{  
  "id": "MyEmrCluster",  
  "type": "EmrCluster",  
  "hadoopVersion": "2.x",  
  "keyPair": "my-key-pair",  
  "masterInstanceType": "m3.xlarge",  
  "coreInstanceType": "m3.large",  
  "coreInstanceCount": "10",  
  "taskInstanceType": "m3.large",  
  "taskInstanceCount": "10",  
  "role": "emrServiceRole",  
  "resourceRole": "emrInstanceProfile"  
}
```

Verwenden Sie die EmrCluster Ressource im AWS SDK for Java

Example

Das folgende Beispiel zeigt, wie Sie mit einem `EmrCluster` und `EmrActivity` einen Amazon EMR 4.x-Cluster erstellen, um einen Spark-Schritt mithilfe des Java-SDK auszuführen:

```
public class dataPipelineEmr4 {  
  
    public static void main(String[] args) {  
  
        AWSCredentials credentials = null;  
        credentials = new ProfileCredentialsProvider("/path/to/  
AwsCredentials.properties", "default").getCredentials();  
    }  
}
```

```
DataPipelineClient dp = new DataPipelineClient(credentials);
CreatePipelineRequest createPipeline = new
CreatePipelineRequest().withName("EMR4SDK").withUniqueId("unique");
CreatePipelineResult createPipelineResult = dp.createPipeline(createPipeline);
String pipelineId = createPipelineResult.getPipelineId();

PipelineObject emrCluster = new PipelineObject()
    .withName("EmrClusterObj")
    .withId("EmrClusterObj")
    .withFields(
        new Field().withKey("releaseLabel").withStringValue("emr-4.1.0"),
        new Field().withKey("coreInstanceCount").withStringValue("3"),
        new Field().withKey("applications").withStringValue("spark"),
        new Field().withKey("applications").withStringValue("Presto-Sandbox"),
        new Field().withKey("type").withStringValue("EmrCluster"),
        new Field().withKey("keyPair").withStringValue("myKeyName"),
        new Field().withKey("masterInstanceType").withStringValue("m3.xlarge"),
        new Field().withKey("coreInstanceType").withStringValue("m3.xlarge")
    );

PipelineObject emrActivity = new PipelineObject()
    .withName("EmrActivityObj")
    .withId("EmrActivityObj")
    .withFields(
        new Field().withKey("step").withStringValue("command-runner.jar,spark-submit,--"
executor-memory,1g,--class,org.apache.spark.examples.SparkPi,/usr/lib/spark/lib/spark-"
examples.jar,10"),
        new Field().withKey("runsOn").withRefValue("EmrClusterObj"),
        new Field().withKey("type").withStringValue("EmrActivity")
    );

PipelineObject schedule = new PipelineObject()
    .withName("Every 15 Minutes")
    .withId("DefaultSchedule")
    .withFields(
        new Field().withKey("type").withStringValue("Schedule"),
        new Field().withKey("period").withStringValue("15 Minutes"),
        new Field().withKey("startAt").withStringValue("FIRST_ACTIVATION_DATE_TIME")
    );

PipelineObject defaultObject = new PipelineObject()
    .withName("Default")
    .withId("Default")
    .withFields(
```

```
new Field().withKey("failureAndRerunMode").withStringValue("CASCADE"),
new Field().withKey("schedule").withRefValue("DefaultSchedule"),
new
Field().withKey("resourceRole").withStringValue("DataPipelineDefaultResourceRole"),
new Field().withKey("role").withStringValue("DataPipelineDefaultRole"),
new Field().withKey("pipelineLogUri").withStringValue("s3://myLogUri"),
new Field().withKey("scheduleType").withStringValue("cron")
);

List<PipelineObject> pipelineObjects = new ArrayList<PipelineObject>();

pipelineObjects.add(emrActivity);
pipelineObjects.add(emrCluster);
pipelineObjects.add(defaultObject);
pipelineObjects.add(schedule);

PutPipelineDefinitionRequest putPipelineDefintion = new PutPipelineDefinitionRequest()
    .withPipelineId(pipelineId)
    .withPipelineObjects(pipelineObjects);

PutPipelineDefinitionResult putPipelineResult =
dp.putPipelineDefinition(putPipelineDefintion);
System.out.println(putPipelineResult);

ActivatePipelineRequest activatePipelineReq = new ActivatePipelineRequest()
    .withPipelineId(pipelineId);
ActivatePipelineResult activatePipelineRes = dp.activatePipeline(activatePipelineReq);

    System.out.println(activatePipelineRes);
    System.out.println(pipelineId);

}

}
```

Einen Amazon EMR-Cluster in einem privaten Subnetz konfigurieren

Example

Dieses Beispiel enthält eine Konfiguration, mit der der Cluster in einem privaten Subnetz in einer VPC gestartet wird. Weitere Informationen finden Sie unter [Starten von Amazon EMR-Clustern in einer VPC](#) im Amazon EMR Management Guide. Diese Konfiguration ist optional. Sie können sie in einer beliebigen Pipeline verwenden, die ein EmrCluster-Objekt nutzt.

Um einen Amazon EMR-Cluster in einem privaten Subnetz zu starten, geben Sie `SubnetId`, `emrManagedMasterSecurityGroupId`, `emrManagedSlaveSecurityGroupId`, und `serviceAccessSecurityGroupId` in Ihrer `EmrCluster` Konfiguration an.

```
{  
    "objects": [  
        {  
            "output": {  
                "ref": "S3BackupLocation"  
            },  
            "input": {  
                "ref": "DDBSourceTable"  
            },  
            "maximumRetries": "2",  
            "name": "TableBackupActivity",  
            "step": "s3://dynamodb-emr-#{myDDBRegion}/emr-ddb-storage-handler/2.1.0/emr-  
ddb-2.1.0.jar,org.apache.hadoop.dynamodb.tools.DynamoDbExport,#{output.directoryPath},#{input.t  
            "id": "TableBackupActivity",  
            "runsOn": {  
                "ref": "EmrClusterForBackup"  
            },  
            "type": "EmrActivity",  
            "resizeClusterBeforeRunning": "false"  
        },  
        {  
            "readThroughputPercent": "#{myDDBReadThroughputRatio}",  
            "name": "DDBSourceTable",  
            "id": "DDBSourceTable",  
            "type": "DynamoDBDataNode",  
            "tableName": "#{myDDBTableName}"  
        },  
        {  
            "directoryPath": "#{myOutputS3Loc}/#{format(@scheduledStartTime, 'YYYY-MM-dd-HH-  
mm-ss')}",  
            "name": "S3BackupLocation",  
            "id": "S3BackupLocation",  
            "type": "S3DataNode"  
        },  
        {  
            "name": "EmrClusterForBackup",  
            "coreInstanceCount": "1",  
            "taskInstanceCount": "1",  
            "taskInstanceType": "m4.xlarge",  
            "subnets": [  
                "SubnetId": "subnet-00000000"  
            ],  
            "securityGroups": [  
                "SecurityGroupId": "sg-00000000"  
            ],  
            "serviceAccessSecurityGroups": [  
                "SecurityGroupId": "sg-00000000"  
            ]  
        }  
    ]  
}
```

```
"coreInstanceType": "m4.xlarge",
"releaseLabel": "emr-4.7.0",
"masterInstanceType": "m4.xlarge",
"id": "EmrClusterForBackup",
"subnetId": "#{mySubnetId}",
"emrManagedMasterSecurityGroupId": "#{myMasterSecurityGroup}",
"emrManagedSlaveSecurityGroupId": "#{mySlaveSecurityGroup}",
"serviceAccessSecurityGroupId": "#{myServiceAccessSecurityGroup}",
"region": "#{myDDBRegion}",
"type": "EmrCluster",
"keyPair": "user-key-pair"
},
{
  "failureAndRerunMode": "CASCADE",
  "resourceRole": "DataPipelineDefaultResourceRole",
  "role": "DataPipelineDefaultRole",
  "pipelineLogUri": "#{myPipelineLogUri}",
  "scheduleType": "ONDEMAND",
  "name": "Default",
  "id": "Default"
}
],
"parameters": [
  {
    "description": "Output S3 folder",
    "id": "myOutputS3Loc",
    "type": "AWS::S3::ObjectKey"
  },
  {
    "description": "Source DynamoDB table name",
    "id": "myDDBTableName",
    "type": "String"
  },
  {
    "default": "0.25",
    "watermark": "Enter value between 0.1-1.0",
    "description": "DynamoDB read throughput ratio",
    "id": "myDDBReadThroughputRatio",
    "type": "Double"
  },
  {
    "default": "us-east-1",
    "watermark": "us-east-1",
    "description": "Region of the DynamoDB table",
  }
]
```

```
        "id": "myDDBRegion",
        "type": "String"
    },
],
"values": {
    "myDDBRegion": "us-east-1",
    "myDDBTableName": "ddb_table",
    "myDDBReadThroughputRatio": "0.25",
    "myOutputS3Loc": "s3://s3_path",
    "mySubnetId": "subnet_id",
    "myServiceAccessSecurityGroup": "service access security group",
    "mySlaveSecurityGroup": "slave security group",
    "myMasterSecurityGroup": "master security group",
    "myPipelineLogUri": "s3://s3_path"
}
}
```

EBS-Volumes zu Cluster-Knoten hinzufügen

Example

Sie können EBS-Volumes an beliebige Knoten im EMR-Cluster innerhalb der Pipeline anfügen. Verwenden Sie zum Anfügen von EBS-Volumes an Knoten `coreEbsConfiguration`, `masterEbsConfiguration` und `taskEbsConfiguration` in Ihrer `EmrCluster`-Konfiguration.

Dieses Beispiel für den Amazon EMR-Cluster verwendet Amazon EBS-Volumes für seine Master-, Task- und Core-Knoten. Weitere Informationen finden Sie unter [Amazon EBS-Volumes in Amazon EMR](#) im Amazon EMR Management Guide.

Diese Konfigurationen sind optional. Sie können sie in beliebigen Pipelines verwenden, die ein `EmrCluster`-Objekt nutzen.

Klicken Sie in der Pipeline auf die `EmrCluster`-Objektkonfiguration und dann auf Master EBS Configuration (Master-EBS-Konfiguration), Core EBS Configuration (Core-EBS-Konfiguration) oder Task EBS Configuration (Aufgaben-EBS-Konfiguration) und geben Sie die Konfigurationsdetails wie im folgenden Beispiel ein.

```
{
    "objects": [
        {
            "output": {
                "ref": "S3BackupLocation"
            }
        }
    ]
}
```

```
},
"input": {
    "ref": "DDBSourceTable"
},
"maximumRetries": "2",
"name": "TableBackupActivity",
"step": "s3://dynamodb-emr-#{myDDBRegion}/emr-ddb-storage-handler/2.1.0/emr-
ddb-2.1.0.jar,org.apache.hadoop.dynamodb.tools.DynamoDbExport,#{output.directoryPath},#{input.t
    "id": "TableBackupActivity",
    "runsOn": {
        "ref": "EmrClusterForBackup"
    },
    "type": "EmrActivity",
    "resizeClusterBeforeRunning": "false"
},
{
    "readThroughputPercent": "#{myDDBReadThroughputRatio}",
    "name": "DDBSourceTable",
    "id": "DDBSourceTable",
    "type": "DynamoDBDataNode",
    "tableName": "#{myDDBTableName}"
},
{
    "directoryPath": "#{myOutputS3Loc}/#{format(@scheduledStartTime, 'YYYY-MM-dd-HH-
mm-ss')}",
    "name": "S3BackupLocation",
    "id": "S3BackupLocation",
    "type": "S3DataNode"
},
{
    "name": "EmrClusterForBackup",
    "coreInstanceCount": "1",
    "taskInstanceCount": "1",
    "taskInstanceType": "m4.xlarge",
    "coreInstanceType": "m4.xlarge",
    "releaseLabel": "emr-4.7.0",
    "masterInstanceType": "m4.xlarge",
    "id": "EmrClusterForBackup",
    "subnetId": "#{mySubnetId}",
    "emrManagedMasterSecurityGroupId": "#{myMasterSecurityGroup}",
    "emrManagedSlaveSecurityGroupId": "#{mySlaveSecurityGroup}",
    "region": "#{myDDBRegion}",
    "type": "EmrCluster",
    "coreEbsConfiguration": {
```

```
        "ref": "EBSConfiguration"
    },
    "masterEbsConfiguration": {
        "ref": "EBSConfiguration"
    },
    "taskEbsConfiguration": {
        "ref": "EBSConfiguration"
    },
    "keyPair": "user-key-pair"
},
{
    "name": "EBSConfiguration",
    "id": "EBSConfiguration",
    "ebsOptimized": "true",
    "ebsBlockDeviceConfig" : [
        { "ref": "EbsBlockDeviceConfig" }
    ],
    "type": "EbsConfiguration"
},
{
    "name": "EbsBlockDeviceConfig",
    "id": "EbsBlockDeviceConfig",
    "type": "EbsBlockDeviceConfig",
    "volumesPerInstance" : "2",
    "volumeSpecification" : {
        "ref": "VolumeSpecification"
    }
},
{
    "name": "VolumeSpecification",
    "id": "VolumeSpecification",
    "type": "VolumeSpecification",
    "sizeInGB": "500",
    "volumeType": "io1",
    "iops": "1000"
},
{
    "failureAndRerunMode": "CASCADE",
    "resourceRole": "DataPipelineDefaultResourceRole",
    "role": "DataPipelineDefaultRole",
    "pipelineLogUri": "#{myPipelineLogUri}",
    "scheduleType": "ONDEMAND",
    "name": "Default",
    "id": "Default"
```

```
        },
    ],
  "parameters": [
    {
      "description": "Output S3 folder",
      "id": "myOutputS3Loc",
      "type": "AWS::S3::ObjectKey"
    },
    {
      "description": "Source DynamoDB table name",
      "id": "myDDBTableName",
      "type": "String"
    },
    {
      "default": "0.25",
      "watermark": "Enter value between 0.1-1.0",
      "description": "DynamoDB read throughput ratio",
      "id": "myDDBReadThroughputRatio",
      "type": "Double"
    },
    {
      "default": "us-east-1",
      "watermark": "us-east-1",
      "description": "Region of the DynamoDB table",
      "id": "myDDBRegion",
      "type": "String"
    }
  ],
  "values": {
    "myDDBRegion": "us-east-1",
    "myDDBTableName": "ddb_table",
    "myDDBReadThroughputRatio": "0.25",
    "myOutputS3Loc": "s3://s3_path",
    "mySubnetId": "subnet_id",
    "mySlaveSecurityGroup": "slave security group",
    "myMasterSecurityGroup": "master security group",
    "myPipelineLogUri": "s3://s3_path"
  }
}
```

Weitere Informationen finden Sie unter:

- [EmrActivity](#)

HttpProxy

HttpProxy ermöglicht es Ihnen, Ihren eigenen Proxy zu konfigurieren und Task Runner über ihn auf den AWS Data Pipeline Dienst zugreifen zu lassen. Es ist nicht erforderlich, einen ausführten TaskRunner mit diesen Informationen zu konfigurieren.

Beispiel für HttpProxy ein TaskRunner

Die folgenden Pipeline-Definition zeigt ein HttpProxy-Objekt:

```
{
  "objects": [
    {
      "schedule": {
        "ref": "Once"
      },
      "pipelineLogUri": "s3://myDPLogUri/path",
      "name": "Default",
      "id": "Default"
    },
    {
      "name": "test_proxy",
      "hostname": "hostname",
      "port": "port",
      "username": "username",
      "*password": "password",
      "windowsDomain": "windowsDomain",
      "type": "HttpProxy",
      "id": "test_proxy",
    },
    {
      "name": "ShellCommand",
      "id": "ShellCommand",
      "runsOn": {
        "ref": "Resource"
      },
      "type": "ShellCommandActivity",
      "command": "echo 'hello world' "
    },
    {
      "period": "1 day",
      "startDateTime": "2013-03-09T00:00:00",
      "name": "Once",
    }
  ]
}
```

```
        "id": "Once",
        "endDateTime": "2013-03-10T00:00:00",
        "type": "Schedule"
    },
{
    "role": "dataPipelineRole",
    "httpProxy": {
        "ref": "test_proxy"
    },
    "actionOnResourceFailure": "retrynone",
    "maximumRetries": "0",
    "type": "Ec2Resource",
    "terminateAfter": "10 minutes",
    "resourceRole": "resourceRole",
    "name": "Resource",
    "actionOnTaskFailure": "terminate",
    "securityGroups": "securityGroups",
    "keyPair": "keyPair",
    "id": "Resource",
    "region": "us-east-1"
}
],
"parameters": []
}
```

Syntax

Pflichtfelder	Description	Slot-Typ
hostname	Der Host des Proxys, über den Clients eine Verbindung zu AWS-Services herstellen.	String
port	Port des Proxy-Hosts, den die Clients verwenden, um eine Verbindung zu AWS-Services herzustellen.	String

Optionale Felder	Description	Slot-Typ
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}
*Passwort	Passwort für den Proxy.	String
s3 NoProxy	Deaktiviert den HTTP-Proxy, wenn eine Verbindung zu Amazon S3 hergestellt wird	Boolesch
username	Benutzername für den Proxy.	String
windowsDomain	Der Windows-Domänenname für NTLM Proxy.	String
windowsWorkgroup	Der Windows-Arbeitsgruppenname für NTLM Proxy.	String

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Vorbedingungen

Im Folgenden sind die AWS Data Pipeline Vorbedingungsobjekte aufgeführt:

Objekte

- [Dynamo DBData ist vorhanden](#)
- [Dynamo ist vorhanden DBTable](#)
- [Vorhanden](#)
- [S3 KeyExists](#)
- [S3 PrefixNotEmpty](#)
- [ShellCommandPrecondition](#)

Dynamo DBData ist vorhanden

Eine Vorbedingung, um zu überprüfen, ob Daten in einer DynamoDB-Tabelle vorhanden sind.

Syntax

Pflichtfelder	Description	Slot-Typ
role	Legt die Rolle für die Ausführung der Vorbedingung fest.	String
tableName	Die zu prüfende DynamoDB-Tabelle.	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum

Optionale Felder	Description	Slot-Typ
failureAndRerunModus	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden	Aufzählung
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlgeschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId"}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId"}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id"}
preconditionTimeout	Der Zeitraum ab dem die Vorbedingung als fehlgeschlagen gekennzeichnet ist, wenn sie noch nicht erfüllt ist	Zeitraum
reportProgressTimeout	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum

Optionale Felder	Description	Slot-Typ
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „activeInstances“: {"ref": "myRunnableObject Id"}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Zuständigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {"ref": "myRunnableObject Id"}
currentRetryCount	Anzahl, wie oft die Vorbedingung in diesem Versuch probiert wurde.	String
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String

Laufzeitfelder	Description	Slot-Typ
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
hostname	Der Hostname des Clients, der den Aufgabenversuch aufnimmt.	String
lastRetryTime	Das letzte Mal, dass die Vorbedingung in diesem Versuch probiert wurde.	String
node	Der Knoten, für den diese Vorbedingung ausgeführt wird	Referenzobjekt, z. B. „node“: {"ref": "myRunnableObject Id"}
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen.	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen.	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „WaitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Dynamo ist vorhanden DBTable

Eine Vorbedingung, um zu überprüfen, ob die DynamoDB-Tabelle existiert.

Syntax

Pflichtfelder	Description	Slot-Typ
role	Legt die Rolle für die Ausführung der Vorbedingung fest.	String
tableName	Die zu prüfende DynamoDB-Tabelle.	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum

Optionale Felder	Description	Slot-Typ
failureAndRerunModus	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden	Aufzählung
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlgeschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId"}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId"}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id"}
preconditionTimeout	Der Zeitraum ab dem die Vorbedingung als fehlgeschlagen gekennzeichnet ist, wenn sie noch nicht erfüllt ist	Zeitraum
reportProgressTimeout	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum

Optionale Felder	Description	Slot-Typ
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „activeInstances“: {"ref": "myRunnableObject Id"}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Zuständigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {"ref": "myRunnableObject Id"}
currentRetryCount	Anzahl, wie oft die Vorbedingung in diesem Versuch probiert wurde.	String
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String

Laufzeitfelder	Description	Slot-Typ
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
hostname	Der Hostname des Clients, der den Aufgabenversuch aufnimmt.	String
lastRetryTime	Das letzte Mal, dass die Vorbedingung in diesem Versuch probiert wurde.	String
node	Der Knoten, für den diese Vorbedingung ausgeführt wird	Referenzobjekt, z. B. „node“: {"ref": "myRunnableObject Id"}
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „WaitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Vorhanden

Prüft, ob eine Datenknotenobjekt vorhanden ist.

 Note

Wir empfehlen, stattdessen die vom System verwalteten Vorbedingungen zu verwenden. Weitere Informationen finden Sie unter [Vorbedingungen](#).

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Das InputData-Objekt verweist auf dieses Objekt, Ready, und auf ein anderes Objekt, das Sie in derselben Pipeline-Definitionsdatei definieren. CopyPeriod ist ein Schedule-Objekt.

```
{
  "id" : "InputData",
  "type" : "S3DataNode",
  "schedule" : { "ref" : "CopyPeriod" },
  "filePath" : "s3://amzn-s3-demo-bucket/InputData/#{@scheduledStartTime.format('YYYY-MM-dd-hh:mm')}.csv",
  "precondition" : { "ref" : "Ready" }
},
{
  "id" : "Ready",
}
```

```

    "type" : "Exists"
}

```

Syntax

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
failureAndRerunMode	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden.	Aufzählung
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlgeschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId"}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId"}

Optionale Felder	Description	Slot-Typ
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}
preconditionTimeout	Der Zeitraum ab dem die Vorbedingung als fehlgeschlagen gekennzeichnet ist, wenn sie noch nicht erfüllt ist	Zeitraum
reportProgressTime out	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „activeInstances“: {"ref": "myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String

Laufzeitfelder	Description	Slot-Typ
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {"ref": "myRunnableObject Id"}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
hostname	Der Hostname des Clients, der den Aufgabenversuch aufnimmt.	String
node	Der Knoten, für den diese Vorbedingung ausgeführt wird.	Referenzobjekt, z. B. „node“: {"ref": "myRunnableObject Id"}
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen.	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen.	DateTime
@Status	Der Status des Objekts.	String

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „WaitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Weitere Informationen finden Sie unter:

- [ShellCommandPrecondition](#)

S3 KeyExists

Prüft, ob ein Schlüssel in einem Amazon S3-Datenknoten vorhanden ist.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Die Vorbedingung wird ausgelöst, wenn der Schlüssel, `s3://amzn-s3-demo-bucket/mykey`, auf den der `s3Key`-Parameter verweist, vorhanden ist.

```
{
```

```

"id" : "InputReady",
"type" : "S3KeyExists",
"role" : "test-role",
"s3Key" : "s3://amzn-s3-demo-bucket/mykey"
}

```

Sie können S3KeyExists auch als Voraussetzung für die zweite Pipeline verwenden, die darauf wartet, dass die erste Pipeline abgeschlossen wird. Gehen Sie hierzu wie folgt vor:

1. Schreiben Sie am Ende der Fertigstellung der ersten Pipeline eine Datei in Amazon S3.
2. Erstellen Sie eine S3KeyExists-Vorbedingung für die zweite Pipeline.

Syntax

Pflichtfelder	Description	Slot-Typ
role	Legt die Rolle für die Ausführung der Vorbedingung fest.	String
s3Key	Der Amazon S3 S3-Schlüssel.	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout vor dem Versuch, die Remote-Arbeit noch einmal auszuführen. Wenn diese Option aktiviert ist, wird erneut versucht, eine Remote-Aktivität durchzuführen, die nach dem Start nicht innerhalb der festgelegten Zeit abgeschlossen wird.	Zeitraum
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlschlagen oder erneut ausgeführt werden.	Aufzählung

Optionale Felder	Description	Slot-Typ
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maximumRetries	Maximale Anzahl der Versuche, die bei einem Fehler initiiert werden.	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId „}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId „}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId „}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id „}
preconditionTimeout	Der Zeitraum ab dem die Vorbedingung als fehlgeschlagen gekennzeichnet ist, wenn sie noch nicht erfüllt ist.	Zeitraum
reportProgressTime out	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in <code>reportProgress</code> . Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei aufeinander folgenden Versuchen.	Zeitraum

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „activeInstances“: {"ref": "myRunnableObject Id"}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn": {"ref": "myRunnableObject Id"}
currentRetryCount	Anzahl, wie oft die Vorbedingung in diesem Versuch probiert wurde.	String
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String

Laufzeitfelder	Description	Slot-Typ
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
lastRetryTime	Das letzte Mal, dass die Vorbedingung in diesem Versuch probiert wurde.	String
node	Der Knoten, für den diese Vorbedingung ausgeführt wird	Referenzobjekt, z. B. „node“: {"ref": "myRunnableObject Id"}
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „WaitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String

Systemfelder	Description	Slot-Typ
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Weitere Informationen finden Sie unter:

- [ShellCommandPrecondition](#)

S3 PrefixNotEmpty

Eine Voraussetzung, um zu überprüfen, ob die Amazon S3 S3-Objekte mit dem angegebenen Präfix (dargestellt als URI) vorhanden sind.

Beispiel

Es folgt ein Beispiel für die Verwendung dieses Objekttyps mit erforderlichen, optionalen und Ausdrucksfeldern.

```
{
  "id" : "InputReady",
  "type" : "S3PrefixNotEmpty",
  "role" : "test-role",
  "s3Prefix" : "#{node.filePath}"
}
```

Syntax

Pflichtfelder	Description	Slot-Typ
role	Legt die Rolle für die Ausführung der Vorbedingung fest.	String
s3Prefix	Das Amazon S3 S3-Präfix zur Überprüfung der Existenz von Objekten.	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum
failureAndRerunMods	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden	Aufzählung
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlgeschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId"}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId"}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id"}

Optionale Felder	Description	Slot-Typ
preconditionTimeout	Der Zeitraum ab dem die Vorbedingung als fehlgeschlagen gekennzeichnet ist, wenn sie noch nicht erfüllt ist	Zeitraum
reportProgressTime out	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum

Laufzeitfelder	Description	Slot-Typ
@activeInstances	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „activeInstances“: {"ref": " myRunnableObject Id "}
@actualEndTime	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {" ref": " myRunnableObject Id "}

Laufzeitfelder	Description	Slot-Typ
currentRetryCount	Anzahl, wie oft die Vorbedingung in diesem Versuch probiert wurde.	String
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
lastRetryTime	Das letzte Mal, dass die Vorbedingung in diesem Versuch probiert wurde.	String
node	Der Knoten, für den diese Vorbedingung ausgeführt wird.	Referenzobjekt, z. B. „node“: {"ref": "myRunnableObject Id"}
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime
@scheduledEndTime	Endzeit für Objekt einplanen.	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen.	DateTime
@Status	Der Status des Objekts.	String

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „WaitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Weitere Informationen finden Sie unter:

- [ShellCommandPrecondition](#)

ShellCommandPrecondition

Ein Unix/Linux Shell-Befehl, der als Vorbedingung ausgeführt werden kann.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
  "id" : "VerifyDataReadiness",
```

```

  "type" : "ShellCommandPrecondition",
  "command" : "perl check-data-ready.pl"
}

```

Syntax

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
command	Den auszuführenden Befehl. Dieser Wert und alle zugehörigen Parameter müssen in der Umgebung funktionieren, in der Sie den Task-Runner ausführen.	String
scriptUri	Ein Amazon S3-URI-Pfad für eine Datei, die heruntergeladen und als Shell-Befehl ausgeführt werden soll. Nur das Feld scriptUri oder das Befehlsfeld sollten vorhanden sein. scriptUri kann keine Parameter verwenden. Verwenden Sie stattdessen das Befehlsfeld.	String

Optionale Felder	Description	Slot-Typ
attemptStatus	Zuletzt gemeldeter Status von der Remote-Aktivität.	String
attemptTimeout	Timeout für die Remote-Arbeit abgeschlossen. Wenn diese Option aktiviert ist, kann eine Remote-Aktivität, die nicht innerhalb der festgelegten Startzeit abgeschlossen wird, wiederholt werden.	Zeitraum

Optionale Felder	Description	Slot-Typ
failureAndRerunModus	Beschreibt das Verhalten des Konsumentenknotens, wenn Abhängigkeiten fehlgeschlagen oder erneut ausgeführt werden	Aufzählung
lateAfterTimeout	Die nach dem Start der Pipeline verstrichene Zeit, innerhalb der das Objekt abgeschlossen werden muss. Sie wird nur ausgelöst, wenn der Zeitplantyp nicht auf eingestellt ist. ondemand	Zeitraum
maximumRetries	Maximale Anzahl von Versuchen bei Ausfällen	Ganzzahl
onFail	Eine Aktion, die ausgeführt werden soll, wenn das aktuelle Objekt fehlgeschlägt.	Referenzobjekt, z. B. „onFail“: {"ref": "myActionId"}
onLateAction	Aktionen, die ausgelöst werden sollen, wenn ein Objekt noch nicht geplant oder noch nicht abgeschlossen wurde.	Referenzobjekt, z. B. "onLateAction": {"ref": "myActionId"}
onSuccess	Eine Aktion, die ausgeführt wird, wenn das aktuelle Objekt erfolgreich ist.	Referenzobjekt, z. B. „onSuccess“: {"ref": "myActionId"}
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id"}
preconditionTimeout	Der Zeitraum ab dem die Vorbedingung als fehlgeschlagen gekennzeichnet ist, wenn sie noch nicht erfüllt ist	Zeitraum
reportProgressTimeout	Timeout für aufeinanderfolgende Aufrufe von Remote-Arbeit in reportProgress. Wenn diese Option aktiviert ist, werden Remote-Aktivitäten, die den Fortschritt für den angegebenen Zeitraum nicht melden, als fehlgeschlagen angesehen und es wird erneut versucht.	Zeitraum

Optionale Felder	Description	Slot-Typ
retryDelay	Die Zeitüberschreitungsdauer zwischen zwei Wiederholungsversuchen.	Zeitraum
scriptArgument	Argument, das an ein Shell-Skript übergeben werden soll	String
stderr	Der Amazon S3 S3-Pfad, der umgeleitete Systemfehlermeldungen vom Befehl empfängt. Wenn Sie das <code>runsOn</code> Feld verwenden, muss es sich um einen Amazon S3 S3-Pfad handeln, da die Ressource, auf der Ihre Aktivität ausgeführt wird, vorübergehend ist. Wenn Sie jedoch das Feld <code>workerGroup</code> angeben, ist ein lokaler Dateipfad zulässig.	String
stdout	Der Amazon S3 S3-Pfad, der die umgeleitete Ausgabe des Befehls empfängt. Wenn Sie das <code>runsOn</code> Feld verwenden, muss es sich um einen Amazon S3 S3-Pfad handeln, da die Ressource, auf der Ihre Aktivität ausgeführt wird, vorübergehend ist. Wenn Sie jedoch das Feld <code>workerGroup</code> angeben, ist ein lokaler Dateipfad zulässig.	String

Laufzeitfelder	Description	Slot-Typ
<code>@activeInstances</code>	Liste der aktuell geplanten aktiven Instance-Objekte.	Referenzobjekt, z. B. „ <code>activeInstances</code> “: {"ref": "Id"} <code>myRunnableObject</code>
<code>@actualEndTime</code>	Zeitpunkt, zu dem die Ausführung dieses Objekts abgeschlossen wurde.	DateTime

Laufzeitfelder	Description	Slot-Typ
@actualStartTime	Zeitpunkt, zu dem die Ausführung dieses Objekts gestartet wurde.	DateTime
cancellationReason	Die cancellationReason, wenn dieses Objekt storniert wurde.	String
@cascadeFailedOn	Beschreibung der Abhängigkeitskette, bei der das Objekt fehlgeschlagen ist.	Referenzobjekt, z. B. "cascadeFailedOn,: {"ref":" myRunnableObject Id "}
emrStepLog	EMR-Schrittprotokolle nur bei EMR-Aktivitätsversuchen verfügbar	String
errorId	Die errorId, wenn dieses Objekt fehlgeschlagen ist.	String
errorMessage	Die errorMessage, wenn dieses Objekt fehlgeschlagen ist.	String
errorStackTrace	Die Fehler-Stack-Ablaufverfolgung., wenn dieses Objekt fehlgeschlagen ist.	String
hadoopJobLog	Hadoop-Jobprotokolle für Versuche für EMR-basierte Aktivitäten verfügbar.	String
hostname	Der Hostname des Clients, der den Aufgabenvorschlag aufnimmt.	String
node	Der Knoten, für den diese Voraussetzung ausgeführt wird	Referenzobjekt, z. B. „node“: {"ref": "myRunnableObject Id "}
reportProgressTime	Der letzte Zeitpunkt, an dem die Remote-Aktivität einen Fortschritt gemeldet hat.	DateTime

Laufzeitfelder	Description	Slot-Typ
@scheduledEndTime	Endzeit für Objekt einplanen	DateTime
@scheduledStartTime	Startzeit für Objekt einplanen	DateTime
@Status	Der Status des Objekts.	String
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String
@waitingOn	Beschreibung der Liste der Abhängigkeiten, auf die dieses Objekt wartet.	Referenzobjekt, z. B. „WaitingOn“: {"ref": "myRunnableObject Id"}

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Weitere Informationen finden Sie unter:

- [ShellCommandActivity](#)
- [Vorhanden](#)

Datenbanken

Im Folgenden sind die AWS Data Pipeline Datenbankobjekte aufgeführt:

Objekte

- [JdbcDatabase](#)
- [RdsDatabase](#)
- [RedshiftDatabase](#)

JdbcDatabase

Definiert eine JDBC-Datenbank.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
  "id" : "MyJdbcDatabase",
  "type" : "JdbcDatabase",
  "connectionString" : "jdbc:redshift://hostname:portnumber/dbname",
  "jdbcDriverClass" : "com.amazon.redshift.jdbc41.Driver",
  "jdbcDriverJarUri" : "s3://redshift-downloads/drivers/RedshiftJDBC41-1.1.6.1006.jar",
  "username" : "user_name",
  "*password" : "my_password"
}
```

Syntax

Pflichtfelder	Description	Slot-Typ
connectionString	Die JDBC-Verbindungszeichenfolge für den Zugriff auf die Datenbank.	String
jdbcDriverClass	Die Treiberklasse, die vor dem Herstellen der JDBC-Verbindung geladen werden soll.	String
*Passwort	Das anzugebende Passwort.	String
username	Der Benutzername, der anzugeben ist, wenn eine Verbindung zur Datenbank hergestellt wird.	String

Optionale Felder	Description	Slot-Typ
databaseName	Name der logischen Datenbank für das Anfügen.	String
jdbcDriverJarUri	Der Amazon S3-Speicherort der JAR-Datei des JDBC-Treibers für die Verbindung mit der Datenbank. AWS Data Pipeline muss über die Leseberechtigung für diese JAR-Datei verfügen.	String
jdbcProperties	Paare der Form A=B, die als Eigenschaften für JDBC-Verbindungen für diese Datenbank festgelegt werden	String
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte	String

Systemfelder	Description	Slot-Typ
	ergeben Instance-Objekte, die Versuchsobjekte ausführen.	

RdsDatabase

Definiert eine Amazon RDS-Datenbank.

Note

RdsDatabase unterstützt Aurora nicht. Verwenden Sie es stattdessen [the section called “JdbcDatabase”](#) für Aurora.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{  
  "id" : "MyRdsDatabase",  
  "type" : "RdsDatabase",  
  "region" : "us-east-1",  
  "username" : "user_name",  
  "*password" : "my_password",  
  "rdsInstanceId" : "my_db_instance_identifier"  
}
```

Für die Oracle-Engine ist das Feld `jdbcDriverJarUri` eine Pflichtangabe. Sie können den folgenden Treiber festlegen: <http://www.oracle.com/technetwork/database/features/jdbc/jdbc-drivers-12c-download-1958347.html>. Für die SQL-Server-Engine ist das Feld `jdbcDriverJarUri` eine Pflichtangabe. Sie können den folgenden Treiber festlegen: <https://www.microsoft.com/en-us/download/details.aspx?displaylang=en&id=11774>. Für die MySQL- und PostgreSQL-Engines ist das Feld `jdbcDriverJarUri` optional.

Syntax

Pflichtfelder	Description	Slot-Typ
*Passwort	Das anzugebende Passwort.	String
rdsInstanceId	Die DBInstanceIdentifier Eigenschaft der DB-Instance.	String
username	Der Benutzername, der anzugeben ist, wenn eine Verbindung zur Datenbank hergestellt wird.	String

Optionale Felder	Description	Slot-Typ
databaseName	Name der logischen Datenbank für das Anfügen.	String
jdbcDriverJarUri	Der Amazon S3-Speicherort der JAR-Datei des JDBC-Treibers für die Verbindung mit der Datenbank. AWS Data Pipeline muss über die Leseberechtigung für diese JAR-Datei verfügen. Für MySQL- und PostgreSQL-Engines wird der Standardtreiber verwendet, wenn dieses Feld nicht angegeben ist. Sie können den Standardwert jedoch mit diesem Feld überschreiben. Für die Oracle- und SQL Server-Engines ist dieses Feld eine Pflichtangabe.	String
jdbcProperties	Paare der Form A=B, die als Eigenschaften für JDBC-Verbindungen für diese Datenbank festgelegt werden	String
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, zum Beispiel „parent“:

Optionale Felder	Description	Slot-Typ
		{"ref": " myBaseObject Id "}
Region	Der Code für die Region, in der die Datenbank vorhanden ist. Beispiel: us-east-1.	String

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

RedshiftDatabase

Definiert eine Amazon Redshift Redshift-Datenbank. RedshiftDatabase stellt die Eigenschaften der Datenbank dar, die von Ihrer Pipeline verwendet wird.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
```

```

"id" : "MyRedshiftDatabase",
"type" : "RedshiftDatabase",
"clusterId" : "myRedshiftClusterId",
"username" : "user_name",
"*password" : "my_password",
"databaseName" : "database_name"
}

```

Standardmäßig nutzt das Objekt den Postgres-Treiber, für den das Feld `clusterId` erforderlich ist. Um den Amazon Redshift Redshift-Treiber zu verwenden, geben Sie stattdessen die Amazon Redshift Redshift-Datenbankverbindungszeichenfolge aus der Amazon Redshift Redshift-Konsole (beginnt mit „`jdbc:redshift:`“) in das Feld ein. `connectionString`

Syntax

Pflichtfelder	Description	Slot-Typ
*Passwort	Das anzugebende Passwort.	String
username	Der Benutzername, der anzugeben ist, wenn eine Verbindung zur Datenbank hergestellt wird.	String

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
clusterId	Die ID, die der Benutzer bei der Erstellung des Amazon Redshift Redshift-Clusters angegeben hat. Wenn der Endpunkt für Ihren Amazon Redshift Redshift-Cluster beispielsweise <code>mydb.example.us-east-1.redshift.amazonaws.com</code> lautet, lautet die korrekte ID. <code>mydb</code> Sie können diesen Wert in der Amazon Redshift-Konsole über "Cluster Identifier" oder "Cluster Name" ermitteln.	String

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
connectionString	Der JDBC-Endpunkt für die Verbindung mit einer Amazon Redshift Redshift-Instance, die einem anderen Konto als der Pipeline gehört. Sie können nicht sowohl connectionString als auch clusterId angeben.	String

Optionale Felder	Description	Slot-Typ
databaseName	Name der logischen Datenbank für das Anfügen.	String
jdbcProperties	Paare der Form A=B müssen als Eigenschaften für JDBC-Verbindungen für diese Datenbank festgelegt werden.	String
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, zum Beispiel „parent“: {"ref": " myBaseObject Id "}
Region	Der Code für die Region, in der die Datenbank vorhanden ist. Beispiel: us-east-1.	Aufzählung

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Datenformate

Im Folgenden sind die AWS Data Pipeline Datenformat-Objekte aufgeführt:

Objekte

- [CSV-Datenformate](#)
- [Custom Data Format](#)
- [Dynamo-Format DBData](#)
- [Dynamo DBExport DataFormat](#)
- [RegEx Datenformat](#)
- [TSV-Datenformate](#)

CSV-Datenformate

Ein durch Kommas getrenntes Datenformat, bei dem das Trennzeichen für Spalten ein Komma und das Datensatztrennzeichen ein Zeilenumbruch ist.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
  "id" : "MyOutputDataType",
  "type" : "CSV",
```

```

"column" : [
    "Name STRING",
    "Score INT",
    "DateOfBirth TIMESTAMP"
]
}

```

Syntax

Optionale Felder	Description	Slot-Typ
column	Spaltenname mit Datentyp, der von jedem Feld für die Daten angegeben wird, die von diesem Datenknoten beschrieben werden. Beispiel: Bei Hostname STRING verwenden Sie für mehrere Werte Spaltennamen und Datentypen, die durch ein Leerzeichen getrennt sind.	String
escapeChar	Ein Zeichen (z. B. "\"), das den Parser anweist, das nächste Zeichen zu ignorieren.	String
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String

Systemfelder	Description	Slot-Typ
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Custom Data Format

Ein benutzerdefiniertes Datenformat, das auf einer Kombination eines bestimmten Spaltentrennzeichens, Datensatztrennzeichens und des Escape-Zeichens basiert.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
  "id" : "MyOutputDataType",
  "type" : "Custom",
  "columnSeparator" : ",",
  "recordSeparator" : "\n",
  "column" : [
    "Name STRING",
    "Score INT",
    "DateOfBirth TIMESTAMP"
  ]
}
```

Syntax

Pflichtfelder	Description	Slot-Typ
columnSeparator	Das Zeichen, mit dem das Ende einer Spalte in einer Datendatei kennzeichnet wird.	String

Optionale Felder	Description	Slot-Typ
column	Spaltenname mit Datentyp, der von jedem Feld für die Daten angegeben wird, die von diesem Datenknoten beschrieben werden. Beispiel: Bei Hostname STRING verwenden Sie für mehrere Werte Spaltennamen und Datentypen, die durch ein Leerzeichen getrennt sind.	String
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}
recordSeparator	Das Zeichen, mit dem das Ende einer Zeile in einer Datendatei kennzeichnet wird, z. B. "\n". Es werden nur einzelne Zeichen unterstützt.	String

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Dynamo-Format DBData

Wendet ein Schema auf eine DynamoDB-Tabelle an, um sie über eine Hive-Abfrage zugänglich zu machen. `DynamoDBDataFormat` wird mit einem `HiveActivity` Objekt und einer `DynamoDBDataNode` Eingabe und Ausgabe verwendet. `DynamoDBDataFormat` erfordert, dass Sie alle Spalten in Ihrer Hive-Abfrage angeben. Mehr Flexibilität bei der Angabe bestimmter Spalten in einer Hive-Abfrage oder Amazon S3 S3-Unterstützung finden Sie unter [Dynamo DBExport DataFormat](#).

 Note

Boolesche DynamoDB-Typen sind nicht zu booleschen Hive-Typen zugeordnet. Es ist aber möglich, DynamoDB-Ganzzahlwerte von 0 oder 1 den booleschen Hive-Typen zuzuordnen.

Beispiel

Das folgende Beispiel zeigt, wie Sie mit `DynamoDBDataFormat` ein Schema einer `DynamoDBDataNode`-Eingabe zuweisen, wodurch ein `HiveActivity`-Objekt auf die Daten nach benannten Spalten zugreifen und die Daten in eine `DynamoDBDataNode`-Ausgabe kopieren kann.

```
{  
  "objects": [  
    {  
      "id" : "Exists.1",  
      "name" : "Exists.1",  
      "type" : "Exists"  
    },  
    {  
      "id" : "DataFormat.1",  
      "name" : "DataFormat.1",  
      "type" : "DynamoDBDataFormat",  
      "column" : [  
        "hash STRING",  
        "range STRING"  
      ]  
    },  
    {  
      "id" : "DynamoDBDataNode.1",  
      "name" : "DynamoDBDataNode.1",  
      "type" : "DynamoDBDataNode",  
      "columns": [  
        "hash STRING",  
        "range STRING"  
      ]  
    }  
  ]  
}
```

```
"tableName" : "$INPUT_TABLE_NAME",
"schedule" : { "ref" : "ResourcePeriod" },
"dataFormat" : { "ref" : "DataFormat.1" }
},
{
  "id" : "DynamoDBDataNode.2",
  "name" : "DynamoDBDataNode.2",
  "type" : "DynamoDBDataNode",
  "tableName" : "$OUTPUT_TABLE_NAME",
  "schedule" : { "ref" : "ResourcePeriod" },
  "dataFormat" : { "ref" : "DataFormat.1" }
},
{
  "id" : "EmrCluster.1",
  "name" : "EmrCluster.1",
  "type" : "EmrCluster",
  "schedule" : { "ref" : "ResourcePeriod" },
  "masterInstanceType" : "m1.small",
  "keyPair" : "$KEYPAIR"
},
{
  "id" : "HiveActivity.1",
  "name" : "HiveActivity.1",
  "type" : "HiveActivity",
  "input" : { "ref" : "DynamoDBDataNode.1" },
  "output" : { "ref" : "DynamoDBDataNode.2" },
  "schedule" : { "ref" : "ResourcePeriod" },
  "runsOn" : { "ref" : "EmrCluster.1" },
  "hiveScript" : "insert overwrite table ${output1} select * from ${input1} ;"
},
{
  "id" : "ResourcePeriod",
  "name" : "ResourcePeriod",
  "type" : "Schedule",
  "period" : "1 day",
  "startDateTime" : "2012-05-04T00:00:00",
  "endDateTime" : "2012-05-05T00:00:00"
}
]
```

Syntax

Optionale Felder	Description	Slot-Typ
column	<p>Der Spaltenname mit dem Datentyp, der von jedem Feld für die Daten angegeben wird, die von diesem Datenknoten beschrieben werden.</p> <p>Beispiel, <code>hostname STRING</code>. Verwenden Sie für mehrere Werte Spaltennamen und Datentypen, die durch ein Leerzeichen getrennt sind.</p>	String
übergeordneter	<p>Das übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.</p>	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}

Laufzeitfelder	Description	Slot-Typ
@Version	Die Pipeline-Version, mit der das Objekt erstellt wurde	String

Systemfelder	Description	Slot-Typ
@error	Der Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Die Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Dynamo DBExport DataFormat

Wendet ein Schema auf eine DynamoDB-Tabelle an, um sie über eine Hive-Abfrage zugänglich zu machen. Verwenden Sie `DynamoDBExportDataFormat` zusammen mit einem `HiveCopyActivity`-Objekt und `DynamoDBDataNode` oder der `S3DataNode`-Ein- und Ausgabe. `DynamoDBExportDataFormat` hat folgende Vorteile:

- Bietet sowohl DynamoDB- als auch Amazon S3 S3-Unterstützung
- Ermöglicht das Filtern von Daten nach bestimmten Spalten in der Hive-Abfrage
- Exportiert alle Attribute aus DynamoDB, auch wenn Sie ein dünnes Schema haben

Note

Boolesche DynamoDB-Typen sind nicht zu booleschen Hive-Typen zugeordnet. Es ist aber möglich, DynamoDB-Ganzzahlwerte von 0 oder 1 den booleschen Hive-Typen zuzuordnen.

Beispiel

Das folgende Beispiel zeigt, wie Sie mit `HiveCopyActivity` und `DynamoDBExportDataFormat` Daten von einem `DynamoDBDataNode` auf einen anderen kopieren können, während gleichzeitig Daten basierend auf einem Zeitstempel gefiltert werden.

```
{  
  "objects": [  
    {  
      "id" : "DataFormat.1",  
      "name" : "DataFormat.1",  
      "type" : "DynamoDBExportDataFormat",  
      "column" : "timeStamp BIGINT"  
    },  
    {  
      "id" : "DataFormat.2",  
      "name" : "DataFormat.2",  
      "type" : "DynamoDBExportDataFormat"  
    },  
    {  
      "id" : "DynamoDBDataNode.1",  
      "name" : "DynamoDBDataNode.1",  
      "type" : "DynamoDBDataNode",  
      "filter": "timeStamp > 1380000000000"  
    }  
  ]  
}
```

```
"tableName" : "item_mapped_table_restore_temp",
"schedule" : { "ref" : "ResourcePeriod" },
"dataFormat" : { "ref" : "DataFormat.1" }
},
{
  "id" : "DynamoDBDataNode.2",
  "name" : "DynamoDBDataNode.2",
  "type" : "DynamoDBDataNode",
  "tableName" : "restore_table",
  "region" : "us_west_1",
  "schedule" : { "ref" : "ResourcePeriod" },
  "dataFormat" : { "ref" : "DataFormat.2" }
},
{
  "id" : "EmrCluster.1",
  "name" : "EmrCluster.1",
  "type" : "EmrCluster",
  "schedule" : { "ref" : "ResourcePeriod" },
  "masterInstanceType" : "m1.xlarge",
  "coreInstanceCount" : "4"
},
{
  "id" : "HiveTransform.1",
  "name" : "Hive Copy Transform.1",
  "type" : "HiveCopyActivity",
  "input" : { "ref" : "DynamoDBDataNode.1" },
  "output" : { "ref" : "DynamoDBDataNode.2" },
  "schedule" : { "ref" : "ResourcePeriod" },
  "runsOn" : { "ref" : "EmrCluster.1" },
  "filterSql" : "`timeStamp` > unix_timestamp(\"#{@scheduledStartTime}\", \"yyyy-MM-dd'T'HH:mm:ss\")"
},
{
  "id" : "ResourcePeriod",
  "name" : "ResourcePeriod",
  "type" : "Schedule",
  "period" : "1 Hour",
  "startDateTime" : "2013-06-04T00:00:00",
  "endDateTime" : "2013-06-04T01:00:00"
}
]
```

Syntax

Optionale Felder	Description	Slot-Typ
column	Spaltenname mit Datentyp, der von jedem Feld für die Daten angegeben wird, die von diesem Datenknoten beschrieben werden. Beispiel: hostname STRING	String
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

RegEx Datenformat

Ein benutzerdefiniertes Datenformat, das durch einen regulären Ausdruck definiert wird.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
  "id" : "MyInputDataType",
  "type" : "RegEx",
  "inputRegEx" : "([^\n]*( [^\n]* ([^\n]* (-|\n|[^\n]*\n)) ([^\n]*|[^\n]*\n) (-|[0-9]*)(-|[0-9]*)(?: ([^\n]*|[^\n]*\n) ([^\n]*|[^\n]*\n))?)",
  "outputFormat" : "%1$s %2$s %3$s %4$s %5$s %6$s %7$s %8$s %9$s",
  "column" : [
    "host STRING",
    "identity STRING",
    "user STRING",
    "time STRING",
    "request STRING",
    "status STRING",
    "size STRING",
    "referer STRING",
    "agent STRING"
  ]
}
```

Syntax

Optionale Felder	Description	Slot-Typ
column	Spaltenname mit Datentyp, der von jedem Feld für die Daten angegeben wird, die von diesem Datenknoten beschrieben werden. Beispiel: Bei Hostname STRING verwenden Sie für mehrere Werte Spaltennamen und Datentypen, die durch ein Leerzeichen getrennt sind.	String
inputRegEx	Der reguläre Ausdruck zum Analysieren einer S3-Eingabedatei. inputRegEx bietet eine Möglichkeit, Spalten aus relativ unstrukturierten Daten in einer Datei abzurufen.	String

Optionale Felder	Description	Slot-Typ
outputFormat	Die Spaltenfelder wurden von %1\$s %2\$s abgerufen inputRegEx, aber mithilfe der Java-Formatierungssyntax als %1\$s %2\$s referenziert.	String
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "Id "} myBaseObject

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

TSV-Datenformate

Ein durch Kommas getrenntes Datenformat, bei dem das Trennzeichen für Spalten ein Tabulatorzeichen und das Datensatztrennzeichen ein Zeilenumbruch ist.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp.

```
{
  "id" : "MyOutputDataType",
  "type" : "TSV",
  "column" : [
    "Name STRING",
    "Score INT",
    "DateOfBirth TIMESTAMP"
  ]
}
```

Syntax

Optionale Felder	Description	Slot-Typ
column	Spaltenname und Datentyp der Daten, die von diesem Datenknoten beschrieben werden. So gibt "Name STRING" eine Spalte mit dem Namen Name und dem Datentyp STRING an. Trennen Sie mehrere Spaltenname/Datentyp-Paare durch Kommas (wie im Beispiel gezeigt).	String
columnSeparator	Das Zeichen, mit dem die Felder einer Spalte von den Feldern der nächsten Spalte getrennt werden. Standardeinstellung: "\t".	String
escapeChar	Ein Zeichen (z. B. "\\"), das den Parser anweist, das nächste Zeichen zu ignorieren.	String
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, zum Beispiel „parent“: {"ref": " myBaseObject Id "}
recordSeparator	Das Zeichen, das die Datensätze voneinander trennt. Standardeinstellung: "\n".	String

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Aktionen

Im Folgenden sind die AWS Data Pipeline Aktionsobjekte aufgeführt:

Objekte

- [SnsAlarm](#)
- [Beenden](#)

SnsAlarm

Sendet eine Amazon SNS SNS-Benachrichtigung, wenn eine Aktivität fehlschlägt oder erfolgreich abgeschlossen wird.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Die Werte für node.input und node.output stammen vom Datenknoten oder der Aktivität, die im Feld onSuccess auf dieses Objekt verweist.

```
{
  "id" : "SuccessNotify",
  "name" : "SuccessNotify",
  "type" : "SnsAlarm",
  "topicArn" : "arn:aws:sns:us-east-1:28619EXAMPLE:ExampleTopic",
  "subject" : "COPY SUCCESS: #{node.@scheduledStartTime}",
  "message" : "Files were copied from #{node.input} to #{node.output}."
}
```

Syntax

Pflichtfelder	Description	Slot-Typ
Nachricht	Der Textkörper der Amazon SNS-Benachrichtigung.	String
role	Die IAM-Rolle für die Erstellung des Amazon SNS-Alarms.	String
subject	Die Betreffzeile der Amazon SNS-Benachrichtigung.	String
topicArn	Der Amazon SNS-Thema-Ziel-ARN für die Nachricht.	String

Optionale Felder	Description	Slot-Typ
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}

Laufzeitfelder	Description	Slot-Typ
node	Der Knoten, für den diese Aktion ausgeführt wird.	Referenzobjekt, z. B. „node“: {"ref": "myNode Id "}

Laufzeitfelder	Description	Slot-Typ
		myRunnableObject Id "}
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Beenden

Eine Aktion, die den Abbruch einer ausstehenden oder noch nicht abgeschlossenen Aktivität, Ressource oder eines Datenknotens auslöst. AWS Data Pipeline versucht, die Aktivität, die Ressource oder den Datenknoten in den Status CANCELLED zu versetzen, wenn nicht mit dem lateAfterTimeout Wert begonnen wird.

Sie können keine Aktionen beenden, die onSuccess-, OnFail- oder onLateAction-Ressourcen beinhalten.

Beispiel

Es folgt ein Beispiel für diesen Objekttyp. Bei diesem Beispiel enthält das Feld onLateAction MyActivity einen Verweis auf die Aktion DefaultAction1. Wenn Sie eine Aktion für onLateAction bereitstellen, müssen Sie auch einen lateAfterTimeout-Wert für den Zeitraum seit dem geplanten Start der Pipeline festlegen, nach dem die Aktivität als verspätet betrachtet wird.

```
{
  "name" : "MyActivity",
  "id" : "DefaultActivity1",
  "schedule" : {
    "ref" : "MySchedule"
  },
  "runsOn" : {
    "ref" : "MyEmrCluster"
  },
  "lateAfterTimeout" : "1 Hours",
  "type" : "EmrActivity",
  "onLateAction" : {
    "ref" : "DefaultAction1"
  },
  "step" : [
    "s3://amzn-s3-demo-bucket/myPath/myStep.jar,firstArg,secondArg",
    "s3://amzn-s3-demo-bucket/myPath/myOtherStep.jar,anotherArg"
  ]
},
{
  "name" : "TerminateTasks",
  "id" : "DefaultAction1",
  "type" : "Terminate"
}
```

Syntax

Optionale Felder	Description	Slot-Typ
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, zum Beispiel „parent“: {"ref": " myBaseObject Id "}

Laufzeitfelder	Description	Slot-Typ
node	Der Knoten, für den diese Aktion ausgeführt wird.	Referenzobjekt, zum Beispiel „node“:

Laufzeitfelder	Description	Slot-Typ
		{"ref": " myRunnableObject Id "}
@Version	Pipeline-Version, mit der das Objekt erstellt wurde	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Plan

Legt den Zeitplan für ein geplantes Ereignis fest, z. B. die Ausführung einer Aktivität.

 Note

Wenn die Startzeit eines Zeitplans in der Vergangenheit liegt, AWS Data Pipeline füllt Ihre Pipeline erneut auf und beginnt sofort mit der Planung von Läufen ab der angegebenen Startzeit. Wählen Sie für Tests/Entwicklung ein relativ kurzes Intervall. Andernfalls wird AWS Data Pipeline versucht, alle Läufe Ihrer Pipeline für dieses Intervall in eine Warteschlange zu stellen und zu planen. AWS Data Pipeline versucht, versehentliche Backfills zu verhindern, wenn die Pipeline-Komponente `scheduledStartTime` vor einem Tag installiert wurde, indem die Pipeline-Aktivierung blockiert wird.

Beispiele

Es folgt ein Beispiel für diesen Objekttyp. Es definiert einen Zeitplan für jede Stunde ab 00:00:00 Uhr am 01.09.2012 bis um 00:00:00 Uhr am 01.10.2012. Der erste Zeitraum endet um 01:00:00 Uhr am 01.09.2012.

```
{  
    "id" : "Hourly",  
    "type" : "Schedule",  
    "period" : "1 hours",  
    "startDateTime" : "2012-09-01T00:00:00",  
    "endDateTime" : "2012-10-01T00:00:00"  
}
```

Die folgende Pipeline startet um FIRST_ACTIVATION_DATE_TIME und wird jede Stunde bis um 22:00:00 Uhr am 25.04.2014 ausgeführt.

```
{  
    "id": "SchedulePeriod",  
    "name": "SchedulePeriod",  
    "startAt": "FIRST_ACTIVATION_DATE_TIME",  
    "period": "1 hours",  
    "type": "Schedule",  
    "endDateTime": "2014-04-25T22:00:00"  
}
```

Die folgende Pipeline startet um FIRST_ACTIVATION_DATE_TIME und wird jede Stunde ausgeführt. Nach dreimaliger Ausführung ist sie abgeschlossen.

```
{  
    "id": "SchedulePeriod",  
    "name": "SchedulePeriod",  
    "startAt": "FIRST_ACTIVATION_DATE_TIME",  
    "period": "1 hours",  
    "type": "Schedule",  
    "occurrences": "3"  
}
```

Die folgende Pipeline beginnt um 22:00:00 Uhr am 25.04.2014, wird stündlich ausgeführt und endet nach dreimaliger Ausführung.

```
{  
    "id": "SchedulePeriod",  
    "name": "SchedulePeriod",  
    "startDateTime": "2014-04-25T22:00:00",  
    "period": "1 hours",  
    "type": "Schedule",  
    "occurrences": "3"  
}
```

On-Demand mit dem Standardobjekt

```
{  
    "name": "Default",  
    "resourceRole": "DataPipelineDefaultResourceRole",  
    "role": "DataPipelineDefaultRole",  
    "scheduleType": "ondemand"  
}
```

On-demand mit explizitem Zeitplanobjekt

```
{  
    "name": "Default",  
    "resourceRole": "DataPipelineDefaultResourceRole",  
    "role": "DataPipelineDefaultRole",  
    "scheduleType": "ondemand"  
,  
{  
    "name": "DefaultSchedule",  
    "type": "Schedule",  
    "id": "DefaultSchedule",  
    "period": "ONDEMAND_PERIOD",  
    "startAt": "ONDEMAND_ACTIVATION_TIME"  
,
```

Die folgenden Beispiele zeigen, wie ein Zeitplan vom Standardobjekt übernommen werden kann, explizit für das Objekt festgelegt werden kann oder durch eine übergeordnete Objektreferenz übergeben werden kann:

Zeitplan vom Standardobjekt übernommen

```
{
```

```
"objects": [
{
  "id": "Default",
  "failureAndRerunMode": "cascade",
  "resourceRole": "DataPipelineDefaultResourceRole",
  "role": "DataPipelineDefaultRole",
  "pipelineLogUri": "s3://myLogsbucket",
  "scheduleType": "cron",
  "schedule": {
    "ref": "DefaultSchedule"
  }
},
{
  "type": "Schedule",
  "id": "DefaultSchedule",
  "occurrences": "1",
  "period": "1 Day",
  "startAt": "FIRST_ACTIVATION_DATE_TIME"
},
{
  "id": "A_Fresh_NewEC2Instance",
  "type": "Ec2Resource",
  "terminateAfter": "1 Hour"
},
{
  "id": "ShellCommandActivity_HelloWorld",
  "runsOn": {
    "ref": "A_Fresh_NewEC2Instance"
  },
  "type": "ShellCommandActivity",
  "command": "echo 'Hello World!''"
}
]
```

Expliziter Zeitplan für das Objekt

```
{
  "objects": [
{
  "id": "Default",
  "failureAndRerunMode": "cascade",
  "resourceRole": "DataPipelineDefaultResourceRole",
```

```
"role": "DataPipelineDefaultRole",
"pipelineLogUri": "s3://myLogsbucket",
"scheduleType": "cron"

},
{
  "type": "Schedule",
  "id": "DefaultSchedule",
  "occurrences": "1",
  "period": "1 Day",
  "startAt": "FIRST_ACTIVATION_DATE_TIME"
},
{
  "id": "A_Fresh_NewEC2Instance",
  "type": "Ec2Resource",
  "terminateAfter": "1 Hour"
},
{
  "id": "ShellCommandActivity_HelloWorld",
  "runsOn": {
    "ref": "A_Fresh_NewEC2Instance"
  },
  "schedule": {
    "ref": "DefaultSchedule"
  },
  "type": "ShellCommandActivity",
  "command": "echo 'Hello World!''"
}
]
```

Zeitplan von übergeordneter Referenz

```
{
  "objects": [
    {
      "id": "Default",
      "failureAndRerunMode": "cascade",
      "resourceRole": "DataPipelineDefaultResourceRole",
      "role": "DataPipelineDefaultRole",
      "pipelineLogUri": "s3://myLogsbucket",
      "scheduleType": "cron"
```

```

},
{
  "id": "parent1",
  "schedule": {
    "ref": "DefaultSchedule"
  }

},
{
  "type": "Schedule",
  "id": "DefaultSchedule",
  "occurrences": "1",
  "period": "1 Day",
  "startAt": "FIRST_ACTIVATION_DATE_TIME"
},
{
  "id": "A_Fresh_NewEC2Instance",
  "type": "Ec2Resource",
  "terminateAfter": "1 Hour"
},
{
  "id": "ShellCommandActivity_HelloWorld",
  "runsOn": {
    "ref": "A_Fresh_NewEC2Instance"
  },
  "parent": {
    "ref": "parent1"
  },
  "type": "ShellCommandActivity",
  "command": "echo 'Hello World!''"
}
]
}

```

Syntax

Pflichtfelder	Description	Slot-Typ
Zeitraum	Die vorgesehene Häufigkeit der Pipeline-Ausführung. Das Format ist "N [Minuten Stunden Tage Wochen Monate]", wobei N eine Zahl gefolgt von einem der Zeitspezi	Zeitraum

Pflichtfelder	Description	Slot-Typ
	fizierer ist. Beispiel: "15 Minuten", führt die Pipeline alle 15 Minuten aus. Der Mindestzeitraum beträgt 15 Minuten und der maximale Zeitraum beträgt 3 Jahre.	

Erforderliche Gruppe (mindestens eine der folgenden ist erforderlich)	Description	Slot-Typ
startAt	Das Datum und der Zeitpunkt, an dem die geplante Pipeline gestartet werden soll. Der gültige Wert ist FIRST_ACTIVATION_DATE_TIME, der zugunsten der Erstellung einer bedarfsgesteuerten Pipeline als veraltet markiert ist.	Aufzählung
startDateTime	Das Datum und die Uhrzeit zum Starten der geplanten Ausführungen. Sie müssen entweder startDateTime oder StartAt verwenden, aber nicht beide.	DateTime

Optionale Felder	Description	Slot-Typ
endDateTime	Das Datum und die Uhrzeit zum Starten der geplanten Ausführungen. Muss ein Datum und eine Uhrzeit nach dem Wert von startDateTime oder StartAt liegen. Das Standardverhalten besteht darin, Ausführungen so lange zu planen, bis die Pipeline heruntergefahren wird.	DateTime

Optionale Felder	Description	Slot-Typ
Ereignisse	Gibt an, wie oft die Pipeline ausgeführt werden soll, nachdem sie aktiviert wurde. Sie können keine Vorkommen mit verwenden. endDateTime	Ganzzahl
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id "}

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@firstActivationTime	Zeit der Objekterstellung.	DateTime
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Dienstprogramme

Die folgenden Dienstprogrammobjekte konfigurieren andere Pipeline-Objekte:

Themen

- [ShellScriptConfig](#)
- [EmrConfiguration](#)
- [Eigenschaft](#)

ShellScriptConfig

Wird zusammen mit einer Aktivität verwendet, um ein Shell-Skript für preActivityTask Config und postActivityTask Config auszuführen. Dieses Objekt ist für [HadoopActivity](#), [HiveActivity](#)[HiveCopyActivity](#), und verfügbar [PigActivity](#). Sie geben einen S3-URI und eine Liste von Argumenten für das Skript an.

Beispiel

A ShellScriptConfig mit Argumenten:

```
{  
    "id" : "ShellScriptConfig_1",  
    "name" : "prescript",  
    "type" : "ShellScriptConfig",  
    "scriptUri": "s3://my-bucket/shell-cleanup.sh",  
    "scriptArgument" : ["arg1","arg2"]  
}
```

Syntax

Dieses Objekt enthält die folgenden Felder.

Optionale Felder	Description	Slot-Typ
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, zum Beispiel „parent“: {"ref": " myBaseObject Id "}
scriptArgument	Eine Liste der Argumente für das Shell-Skript	String

Optionale Felder	Description	Slot-Typ
scriptUri	Der URI des Skripts in Amazon S3, das heruntergeladen und ausgeführt werden soll.	String

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

EmrConfiguration

Das EmrConfiguration Objekt ist die Konfiguration, die für EMR-Cluster mit Versionen 4.0.0 oder höher verwendet wird. Konfigurationen (als Liste) sind ein Parameter für den RunJobFlow API-Aufruf. Die Konfigurations-API für Amazon EMR verwendet eine Klassifizierung und Eigenschaften. AWS Data Pipeline verwendet EmrConfiguration mit entsprechenden Property-Objekten, um eine [EmrCluster](#) Anwendung wie Hadoop, Hive, Spark oder Pig auf EMR-Clustern zu konfigurieren, die in einer Pipeline-Ausführung gestartet wurden. Da die Konfiguration nur für neue Cluster geändert werden kann, können Sie kein EmrConfiguration Objekt für vorhandene Ressourcen bereitstellen. Weitere Informationen finden Sie unter <https://docs.aws.amazon.com/ElasticMapReduce/latest/ReleaseGuide/>.

Beispiel

Das folgende Konfigurationsobjekt legt die `io.file.buffer.size` und `fs.s3.block.size` Eigenschaften in `core-site.xml` fest:

```
[  
  {  
    "classification": "core-site",  
    "properties":  
    {  
      "io.file.buffer.size": "4096",  
      "fs.s3.block.size": "67108864"  
    }  
  }  
]
```

Die entsprechende Pipeline-Objektdefinition verwendet ein EmrConfiguration Objekt und eine Liste von Property-Objekten im `property` Feld:

```
{  
  "objects": [  
    {  
      "name": "ReleaseLabelCluster",  
      "releaseLabel": "emr-4.1.0",  
      "applications": ["spark", "hive", "pig"],  
      "id": "ResourceId_I1mCc",  
      "type": "EmrCluster",  
      "configuration": {  
        "ref": "coresite"  
      }  
    },  
    {  
      "name": "coresite",  
      "id": "coresite",  
      "type": "EmrConfiguration",  
      "classification": "core-site",  
      "property": [{  
        "ref": "io-file-buffer-size"  
      },  
      {  
        "ref": "fs-s3-block-size"  
      }  
    ]  
  ]
```

```
},
{
  "name": "io-file-buffer-size",
  "id": "io-file-buffer-size",
  "type": "Property",
  "key": "io.file.buffer.size",
  "value": "4096"
},
{
  "name": "fs-s3-block-size",
  "id": "fs-s3-block-size",
  "type": "Property",
  "key": "fs.s3.block.size",
  "value": "67108864"
}
]
```

Im folgenden Beispiel wird eine verschachtelte Konfiguration verwendet, um die Hadoop-Umgebung mit der `hadoop-env`-Klassifizierung festzulegen:

```
[
{
  "classification": "hadoop-env",
  "properties": {},
  "configurations": [
    {
      "classification": "export",
      "properties": {
        "YARN_PROXYSERVER_HEAPSIZE": "2396"
      }
    }
  ]
}
```

Nachfolgend ist das entsprechende Pipeline-Definitionsobjekt mit dieser Konfiguration:

```
{
  "objects": [
    {
      "name": "ReleaseLabelCluster",
      "releaseLabel": "emr-4.0.0",
      "configuration": {
        "classification": "hadoop-env",
        "properties": {},
        "configurations": [
          {
            "classification": "export",
            "properties": {
              "YARN_PROXYSERVER_HEAPSIZE": "2396"
            }
          }
        ]
      }
    }
  ]
}
```

```
"applications": ["spark", "hive", "pig"],
"id": "ResourceId_I1mCc",
"type": "EmrCluster",
"configuration": {
    "ref": "hadoop-env"
},
{
    "name": "hadoop-env",
    "id": "hadoop-env",
    "type": "EmrConfiguration",
    "classification": "hadoop-env",
    "configuration": {
        "ref": "export"
    }
},
{
    "name": "export",
    "id": "export",
    "type": "EmrConfiguration",
    "classification": "export",
    "property": {
        "ref": "yarn-proxyserver-heapsize"
    }
},
{
    "name": "yarn-proxyserver-heapsize",
    "id": "yarn-proxyserver-heapsize",
    "type": "Property",
    "key": "YARN_PROXYSERVER_HEAPSIZE",
    "value": "2396"
},
]
}
```

Im folgenden Beispiel wird eine HIVE-spezifische Eigenschaft für einen EMR-Cluster geändert:

```
{
    "objects": [
        {
            "name": "hivesite",
            "id": "hivesite",
            "type": "EmrConfiguration",
            "classification": "hive-site"
        }
    ]
}
```

```

"classification": "hive-site",
"property": [
    {
        "ref": "hive-client-timeout"
    }
],
{
    "name": "hive-client-timeout",
    "id": "hive-client-timeout",
    "type": "Property",
    "key": "hive.metastore.client.socket.timeout",
    "value": "2400s"
}
]
}

```

Syntax

Dieses Objekt enthält die folgenden Felder.

Pflichtfelder	Description	Slot-Typ
Klassifizierung	Klassifizierung für die Konfiguration.	String

Optionale Felder	Description	Slot-Typ
Konfiguration	Unterkonfiguration für diese Konfiguration.	Referenzobjekt, z. B. „configuration“: {"ref": "Id"} myEmrConfiguration
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, z. B. „parent“: {"ref": "myBaseObject Id"}

Optionale Felder	Description	Slot-Typ
property	Konfigurationseigenschaft	Referenzobjekt, z. B. „Eigenschaft“: {"ref": "myPropertyId"}

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde.	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts	String
@pipelineId	Id der Pipeline, zu der dieses Objekt gehört	String
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen	String

Weitere Informationen finden Sie unter:

- [EmrCluster](#)
- [Eigenschaft](#)
- [Amazon EMR-Versionshinweise](#)

Eigenschaft

Eine einzelne Schlüssel-Wert-Eigenschaft zur Verwendung mit einem Objekt EmrConfiguration .

Beispiel

Die folgende Pipeline-Definition zeigt ein EmrConfiguration Objekt und die entsprechenden Eigenschaftsobjekte zum Starten eines: EmrCluster

```
{  
  "objects": [  
    {  
      "name": "ReleaseLabelCluster",  
      "releaseLabel": "emr-4.1.0",  
      "applications": ["spark", "hive", "pig"],  
      "id": "ResourceId_I1mCc",  
      "type": "EmrCluster",  
      "configuration": {  
        "ref": "coresite"  
      }  
    },  
    {  
      "name": "coresite",  
      "id": "coresite",  
      "type": "EmrConfiguration",  
      "classification": "core-site",  
      "property": [{  
        "ref": "io-file-buffer-size"  
      },  
      {  
        "ref": "fs-s3-block-size"  
      }  
    ],  
    {  
      "name": "io-file-buffer-size",  
      "id": "io-file-buffer-size",  
      "type": "Property",  
      "key": "io.file.buffer.size",  
      "value": "4096"  
    },  
    {  
      "name": "fs-s3-block-size",  
      "id": "fs-s3-block-size",  
      "type": "Property",  
      "key": "fs.s3.block.size",  
      "value": "67108864"  
    }  
  ]  
}
```

```

    }
]
}

```

Syntax

Dieses Objekt enthält die folgenden Felder.

Pflichtfelder	Description	Slot-Typ
Schlüssel	Schlüssel	Zeichenfolge
value	value	String

Optionale Felder	Description	Slot-Typ
übergeordneter	Übergeordnetes Objekt des aktuellen Objekts, aus dem Slots übernommen werden.	Referenzobjekt, zum Beispiel „parent“: {"ref": " myBaseObject Id "}

Laufzeitfelder	Description	Slot-Typ
@Version	Pipeline-Version, mit der das Objekt erstellt wurde	String

Systemfelder	Description	Slot-Typ
@error	Fehler mit einer Beschreibung des falsch formatierten Objekts.	String
@pipelineId	ID der Pipeline, zu der dieses Objekt gehört.	String

Systemfelder	Description	Slot-Typ
@sphere	Die Kugel eines Objekts bezeichnet seinen Platz im Lebenszyklus: Komponentenobjekte ergeben Instance-Objekte, die Versuchsobjekte ausführen.	String

Weitere Informationen finden Sie unter:

- [EmrCluster](#)
- [EmrConfiguration](#)
- [Amazon EMR-Versionshinweise](#)

Arbeiten mit Task Runner

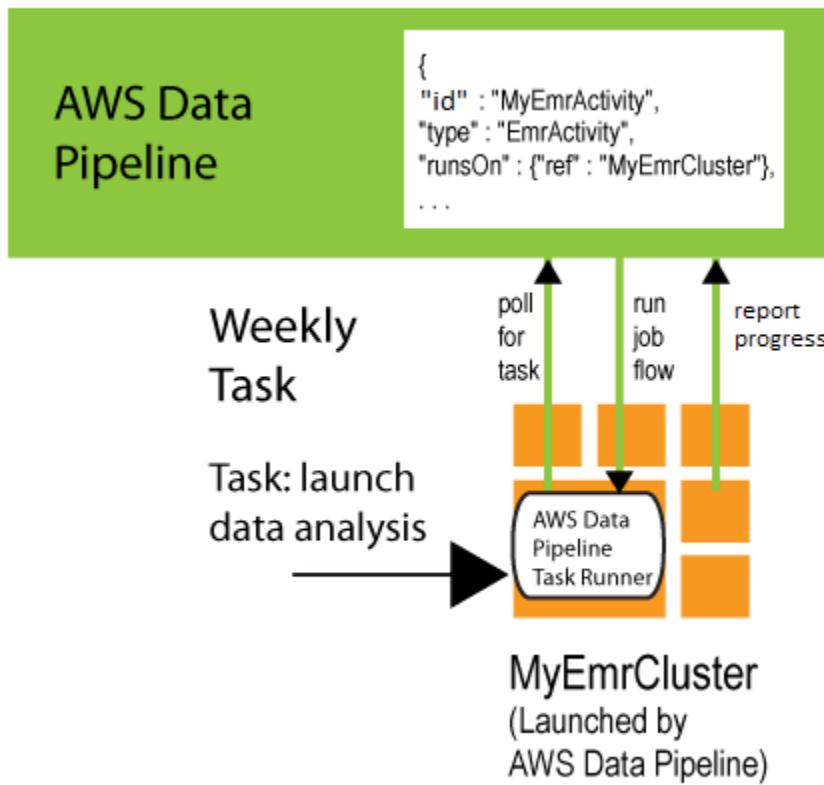
Task Runner ist eine Task-Agent-Anwendung, die geplante Aufgaben abfragt AWS Data Pipeline und sie auf EC2 Amazon-Instances, Amazon EMR-Clustern oder anderen Rechenressourcen ausführt und dabei den Status meldet. Je nach Anwendung können Sie:

- Erlauben AWS Data Pipeline Sie, eine oder mehrere Task Runner-Anwendungen für Sie zu installieren und zu verwalten. Wenn eine Pipeline aktiviert ist, wird automatisch der Standard `Ec2Instance` oder das `EmrCluster` Objekt erstellt, auf das ein `RunsOn`-Feld der Aktivität verweist. AWS Data Pipeline kümmert sich um die Installation von Task Runner auf einer EC2 Instanz oder auf dem Master-Knoten eines EMR-Clusters. In diesem Muster AWS Data Pipeline kann der Großteil der Instanz- oder Clusterverwaltung für Sie erledigt werden.
- Führen Sie alle oder Teile einer Pipeline für von Ihnen verwaltete Ressourcen aus. Zu den potenziellen Ressourcen gehören eine lang laufende EC2 Amazon-Instance, ein Amazon EMR-Cluster oder ein physischer Server. Sie können einen Task-Runner (der entweder Task Runner oder ein benutzerdefinierter Task-Agent Ihres eigenen Geräts sein kann) fast überall installieren, vorausgesetzt, er kann mit dem AWS Data Pipeline Webservice kommunizieren. In diesem Muster übernehmen Sie fast die vollständige Kontrolle darüber, welche Ressourcen verwendet und wie sie verwaltet werden, und Sie müssen Task Runner manuell installieren und konfigurieren. Verwenden Sie dazu die Verfahren in diesem Abschnitt, wie in [Arbeiten an vorhandenen Ressourcen mithilfe von Task Runner ausführen](#) beschrieben.

Task Runner auf AWS Data Pipeline-Verwaltete Ressourcen

Wenn eine Ressource von gestartet und verwaltet wird AWS Data Pipeline, installiert der Webdienst automatisch Task Runner auf dieser Ressource, um Aufgaben in der Pipeline zu verarbeiten.

Sie geben eine Rechenressource (entweder eine EC2 Amazon-Instance oder ein Amazon EMR-Cluster) für das `runsOn` Feld eines Aktivitätsobjekts an. Wenn diese Ressource AWS Data Pipeline gestartet wird, installiert sie Task Runner auf dieser Ressource und konfiguriert sie so, dass alle Aktivitätsobjekte verarbeitet werden, deren `runsOn` Feld auf diese Ressource gesetzt ist. Wenn die Ressource AWS Data Pipeline beendet wird, werden die Task Runner-Protokolle an einem Amazon S3 S3-Speicherort veröffentlicht, bevor sie heruntergefahren wird.



Wenn Sie beispielsweise in einer Pipeline die `EmrActivity` verwenden und im Feld `runsOn` eine `EmrCluster`-Ressource angeben. Bei der AWS Data Pipeline Verarbeitung dieser Aktivität wird ein Amazon EMR-Cluster gestartet und Task Runner auf dem Master-Knoten installiert. Dieser Task Runner verarbeitet dann die Aufgaben für Aktivitäten, deren `runsOn` Feld auf dieses `EmrCluster` Objekt festgelegt ist. Der folgende Ausschnitt aus einer Pipeline-Definition zeigt diese Beziehung zwischen den beiden Objekten.

```
{
  "id" : "MyEmrActivity",
  "name" : "Work to perform on my data",
  "type" : "EmrActivity",
  "runsOn" : {"ref" : "MyEmrCluster"},
  "preStepCommand" : "scp remoteFiles localFiles",
  "step" : "s3://amzn-s3-demo-bucket/myPath/myStep.jar,firstArg,secondArg",
  "step" : "s3://amzn-s3-demo-bucket/myPath/myOtherStep.jar,anotherArg",
  "postStepCommand" : "scp localFiles remoteFiles",
  "input" : {"ref" : "MyS3Input"},
  "output" : {"ref" : "MyS3Output"}
},
{
```

```
"id" : "MyEmrCluster",
"name" : "EMR cluster to perform the work",
"type" : "EmrCluster",
"hadoopVersion" : "0.20",
"keypair" : "myKeyPair",
"masterInstanceType" : "m1.xlarge",
"coreInstanceType" : "m1.small",
"coreInstanceCount" : "10",
"taskInstanceType" : "m1.small",
"taskInstanceCount": "10",
"bootstrapAction" : "s3://elasticmapreduce/libs/ba/configure-hadoop,arg1,arg2,arg3",
"bootstrapAction" : "s3://elasticmapreduce/libs/ba/configure-other-stuff,arg1,arg2"
}
```

Informationen und Beispiele für die Ausführung dieser Aktivitäten finden Sie unter [EmrActivity](#).

Wenn Sie mehrere AWS Data Pipeline verwaltete Ressourcen in einer Pipeline haben, wird Task Runner auf jeder von ihnen installiert, und alle Ressourcen fragen nach zu AWS Data Pipeline verarbeitenden Aufgaben ab.

Arbeiten an vorhandenen Ressourcen mithilfe von Task Runner ausführen

Sie können Task Runner auf Rechenressourcen installieren, die Sie verwalten, z. B. einer EC2 Amazon-Instance oder einem physischen Server oder einer Workstation. Task Runner kann überall auf jeder kompatiblen Hardware oder jedem kompatiblen Betriebssystem installiert werden, vorausgesetzt, es kann mit dem AWS Data Pipeline Webservice kommunizieren.

Dieser Ansatz kann nützlich sein, wenn Sie beispielsweise Daten verarbeiten AWS Data Pipeline möchten, die in der Firewall Ihres Unternehmens gespeichert sind. Durch die Installation von Task Runner auf einem Server im lokalen Netzwerk können Sie sicher auf die lokale Datenbank zugreifen und dann die AWS Data Pipeline nächste auszuführende Aufgabe abfragen. Wenn die Verarbeitung AWS Data Pipeline beendet oder die Pipeline gelöscht wird, läuft die Task Runner-Instanz weiterhin auf Ihrer Rechenressource, bis Sie sie manuell herunterfahren. Die Task Runner-Protokolle bleiben bestehen, nachdem die Pipeline-Ausführung abgeschlossen ist.

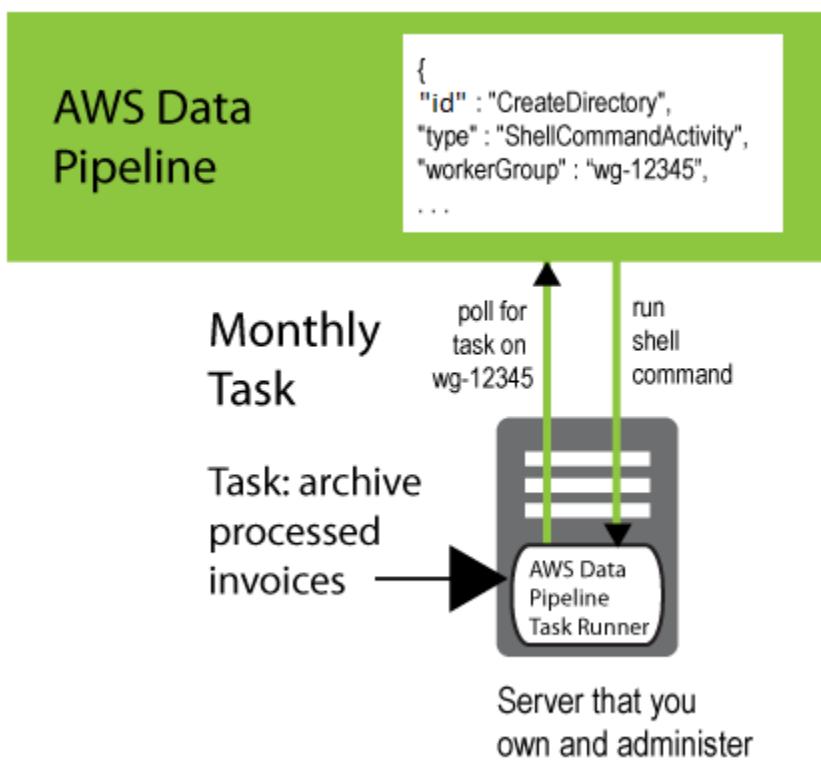
Um Task Runner auf einer von Ihnen verwalteten Ressource zu verwenden, müssen Sie zuerst Task Runner herunterladen und ihn dann mithilfe der Verfahren in diesem Abschnitt auf Ihrer Rechenressource installieren.

Note

Sie können Task Runner nur unter Linux, UNIX oder macOS installieren. Task Runner wird auf dem Windows-Betriebssystem nicht unterstützt.

Um Task Runner 2.0 verwenden zu können, ist mindestens die Java-Version 1.7 erforderlich.

Um einen Task Runner, den Sie installiert haben, mit den Pipeline-Aktivitäten zu verbinden, die er verarbeiten soll, fügen Sie dem Objekt ein `workerGroup` Feld hinzu und konfigurieren Sie Task Runner so, dass er diesen Workergruppenwert abfragt. Dazu übergeben Sie die Worker-Group-Zeichenfolge als Parameter (z. B. --`workerGroup=wg-12345`), wenn Sie die Task Runner-JAR-Datei ausführen.



```
{  
  "id" : ".CreateDirectory",  
  "type" : "ShellCommandActivity",  
  "workerGroup" : "wg-12345",  
  "command" : "mkdir new-directory"  
}
```

Task Runner installieren

In diesem Abschnitt wird erklärt, wie Task Runner installiert und konfiguriert wird und welche Voraussetzungen dafür erforderlich sind. Die Installation ist ein einfacher manueller Prozess.

Um Task Runner zu installieren

1. Task Runner benötigt die Java-Versionen 1.6 oder 1.8. Verwenden Sie den folgenden Befehl, um festzustellen, ob Java installiert ist und welche Version ausgeführt wird:

```
java -version
```

Wenn Sie Java 1.6 oder 1.8 nicht auf Ihrem Computer installiert haben, laden Sie eine dieser Versionen von <http://www.oracle.com/technetwork/java/index.html>. Laden Sie Java herunter und installieren Sie es und fahren Sie mit dem nächsten Schritt fort.

2. Laden Sie es TaskRunner-1.0.jar von <https://s3.amazonaws.com/datapipeline-us-east-1/us-east-1/software/latest/TaskRunner/TaskRunner-1.0.jar> herunter und kopieren Sie es dann in einen Ordner auf der Ziel-Computerressource. Für Amazon EMR-Cluster, die EmrActivity Aufgaben ausführen, installieren Sie Task Runner auf dem Master-Knoten des Clusters.
3. Wenn Sie Task Runner verwenden, um eine Verbindung zum AWS Data Pipeline Webservice herzustellen, um Ihre Befehle zu verarbeiten, benötigen Benutzer programmatischen Zugriff auf eine Rolle, die über Berechtigungen zum Erstellen oder Verwalten von Daten-Pipelines verfügt. Weitere Informationen finden Sie unter [Erteilen programmgesteuerten Zugriffs](#).
4. Task Runner stellt über HTTPS eine Verbindung zum AWS Data Pipeline Webdienst her. Wenn Sie eine AWS Ressource verwenden, stellen Sie sicher, dass HTTPS in der entsprechenden Routingtabelle und Subnetz-ACL aktiviert ist. Wenn Sie eine Firewall oder einen Proxy verwenden, stellen Sie sicher, dass der Port 443 geöffnet ist.

(Optional) Task Runner-Zugriff auf Amazon RDS gewähren

Mit Amazon RDS können Sie den Zugriff auf Ihre DB-Instances mithilfe von Datenbanksicherheitsgruppen (DB-Sicherheitsgruppen) steuern. DB-Sicherheitsgruppen funktionieren wie eine Firewall. Sie steuern den Netzwerkzugriff auf die DB-Instance. Der Netzwerkzugriff auf Ihre DB-Instances ist standardmäßig deaktiviert. Sie müssen Ihre DB-

Sicherheitsgruppen ändern, damit Task Runner auf Ihre Amazon RDS-Instances zugreifen kann. Task Runner erhält Amazon RDS-Zugriff von der Instance, auf der es ausgeführt wird. Daher hängen die Konten und Sicherheitsgruppen, die Sie zu Ihrer Amazon RDS-Instance hinzufügen, davon ab, wo Sie Task Runner installieren.

Um Zugriff auf Task Runner in EC2 -Classic zu gewähren

1. Öffnen Sie die Amazon RDS-Konsole.
2. Klicken Sie im Navigationsbereich auf Instances (DB-Instances) und wählen Sie anschließend Ihre DB-Instance aus.
3. Wählen Sie unter Sicherheit und Netzwerk die Sicherheitsgruppe aus, die die Seite Sicherheitsgruppen mit dieser ausgewählten DB-Sicherheitsgruppe öffnet. Klicken Sie auf das Detailsymbol für die DB-Sicherheitsgruppe.
4. Legen Sie unter Sicherheitsgruppendetails eine Regel mit dem entsprechenden Verbindungstyp und Details an. Diese Felder hängen davon ab, wo Task Runner ausgeführt wird, wie hier beschrieben:
 - Ec2Resource
 - Verbindungstyp: EC2 Security Group
Details: *my-security-group-name* (der Name der Sicherheitsgruppe, die Sie für die EC2 Instanz erstellt haben)
 - EmrResource
 - Verbindungstyp: EC2 Security Group
Details: ElasticMapReduce-master
 - Verbindungstyp: EC2 Security Group
 - Details: ElasticMapReduce-slave
 - Ihre lokale Umgebung (lokal)
 - Verbindungstyp: CIDR/IP:
Details: *my-ip-address* (die IP-Adresse Ihres Computers oder der IP-Adressbereich Ihres Netzwerks, falls sich Ihr Computer hinter einer Firewall befindet)
5. Klicken Sie auf Add (Hinzufügen).

Um Zugriff auf Task Runner in VPC EC2 zu gewähren

1. Öffnen Sie die Amazon RDS-Konsole.
2. Wählen Sie im Navigationsbereich Instances aus.
3. Klicken Sie auf das Detailsymbol für die DB-Sicherheitsgruppe. Öffnen Sie unter Sicherheit und Netzwerk den Link zur Sicherheitsgruppe, über den Sie zur EC2 Amazon-Konsole gelangen. Wenn Sie das alte Konsolendesign für Sicherheitsgruppen verwenden, wechseln Sie zum neuen Konsolendesign, indem Sie auf das Symbol klicken, das oben auf der Konsolenseite angezeigt wird.
4. Wählen Sie auf der Registerkarte Inbound die Option Edit, Add Rule. Geben Sie den Datenbankanschluss an, den Sie beim Starten der DB-Instance verwendet haben. Die Quelle hängt davon ab, wo Task Runner ausgeführt wird, wie hier beschrieben:
 - Ec2Resource
 - *my-security-group-id*(die ID der Sicherheitsgruppe, die Sie für die EC2 Instanz erstellt haben)
 - EmrResource
 - *master-security-group-id*(die ID der ElasticMapReduce-master Sicherheitsgruppe)
 - *slave-security-group-id*(die ID der ElasticMapReduce-slave Sicherheitsgruppe)
 - Ihre lokale Umgebung (lokal)
 - *ip-address*(die IP-Adresse Ihres Computers oder der IP-Adressbereich Ihres Netzwerks, falls sich Ihr Computer hinter einer Firewall befindet)
5. Klicken Sie auf Speichern.

Task Runner starten

Starten Sie Task Runner in einem neuen Befehlszeilenfenster, das sich auf das Verzeichnis bezieht, in dem Sie Task Runner installiert haben, mit dem folgenden Befehl.

```
java -jar TaskRunner-1.0.jar --config ~/credentials.json --workerGroup=myWorkerGroup --region=MyRegion --logUri=s3://amzn-s3-demo-bucket/foldername
```

Die Option --config zeigt auf Ihre Anmeldedaten-Datei.

Die Option `--workerGroup` gibt den Namen Ihrer Auftragnehmergruppe an, der für die zu verarbeitenden Aufgaben in Ihrer Pipeline den gleichen Wert haben muss.

Die Option `--region` gibt den Servicebereich an, von dem aus Aufgaben ausgeführt werden sollen.

Die `--logUri` Option wird verwendet, um Ihre komprimierten Protokolle an einen Speicherort in Amazon S3 zu übertragen.

Wenn Task Runner aktiv ist, gibt es den Pfad, in den die Protokolldateien geschrieben werden, im Terminalfenster aus. Im Folgenden wird ein Beispiel gezeigt.

```
Logging to /Computer_Name/.../output/logs
```

Task-Runner sollte von Ihrer Anmelde-Shell getrennt ausgeführt werden. Wenn Sie eine Terminalanwendung verwenden, um eine Verbindung zu Ihrem Computer herzustellen, müssen Sie möglicherweise ein Dienstprogramm wie nohup oder screen verwenden, um zu verhindern, dass die Task-Runner-Anwendung beendet wird, wenn Sie sich abmelden. Weitere Informationen zu diesen Befehlszeilenoptionen finden Sie unter [Task Runner-Konfigurationsoptionen](#).

Die Task Runner-Protokollierung wird überprüft

Der einfachste Weg, um zu überprüfen, ob Task Runner funktioniert, besteht darin, zu überprüfen, ob er Protokolldateien schreibt. Task Runner schreibt stündlich Protokolldateien in das Verzeichnis `output/logs`, unter dem Verzeichnis, in dem Task Runner installiert ist. Der Dateiname ist `Task_Runner.log.YYYY-MM-DD-HH`, wobei HH von 00 bis 23 in UDT läuft. Um Speicherplatz zu sparen, werden alle Protokolldateien, die älter als acht Stunden sind, mit komprimiert GZip.

Task Runner-Threads und Vorbedingungen

Task Runner verwendet für jede Aufgabe, Aktivität und Vorbedingung einen Threadpool. Die Standardeinstellung für `--tasks` ist 2, was bedeutet, dass zwei Threads aus dem Aufgabenpool zugewiesen sind und jeder Thread den AWS Data Pipeline Dienst nach neuen Aufgaben abfragt. Daher ist `--tasks` ein Attribut für die Leistungsoptimierung, mit dem der Pipeline-Durchsatz optimiert werden kann.

Die Pipeline-Wiederholungslogik für Vorbedingungen erfolgt in Task Runner. Zwei Vorbedingungs-Threads sind AWS Data Pipeline für die Abfrage von Vorbedingungsobjekten zugewiesen. Task Runner berücksichtigt die Felder `RetryDelay` und `PreconditionTimeout` der Vorbedingungsobjekte, die Sie unter Vorbedingungen definieren.

In vielen Fällen kann das Reduzieren des Timeouts für Abfragen und die Anzahl der Wiederholungen dazu beitragen, die Leistung Ihrer Anwendung zu verbessern. In ähnlicher Weise müssen Anwendungen mit Langzeitvorbedingungen die Werte für Timeout und Wiederholung erhöhen. Weitere Informationen zu Vorbedingungsobjekten finden Sie unter [Vorbedingungen](#).

Task Runner-Konfigurationsoptionen

Dies sind die Konfigurationsoptionen, die in der Befehlszeile verfügbar sind, wenn Sie Task Runner starten.

Befehlszeilen-Parameter	Beschreibung
--help	Befehlszeilenhilfe. Beispiel: <code>Java -jar TaskRunner-1.0.jar --help</code>
--config	Pfad und Dateiname Ihrer Datei <code>credentials.json</code> .
--accessId	Ihre AWS Zugriffsschlüssel-ID für Task Runner, die Sie bei Anfragen verwenden können.
--secretKey	Die <code>--secretKey</code> Optionen <code>--accessID</code> und bieten eine Alternative zur Verwendung einer <code>Credentials.json</code> -Datei. Wenn auch eine <code>credentials.json</code> -Datei vorgesehen ist, haben die Optionen <code>--accessID</code> und <code>--secretKey</code> Vorrang.
--endpoint	Ein Endpunkt ist eine URL, die als Eintrittspunkt für einen Webservice fungiert. Der AWS Data Pipeline Service-Endpunkt in der Region, in der Sie Anfragen stellen. Optional. Im Allgemein en reicht es aus, eine Region anzugeben, und Sie müssen den Endpunkt nicht festlegen. Eine

Befehlszeilen-Parameter	Beschreibung
	<p>Liste der AWS Data Pipeline Regionen und Endpunkte finden Sie unter AWS Data Pipeline Regions and Endpoints in der Allgemeine AWS-Referenz</p>
--workerGroup	<p>Der Name der Auftragnehmergruppe, für die Task Runner Aufträge abruft Erforderlich</p> <p>Wenn Task Runner den Webdienst abfragt, verwendet er die von Ihnen angegebenen Anmeldeinformationen und den Wert von, <code>workerGroup</code> um auszuwählen, welche (falls vorhanden) Aufgaben abgerufen werden sollen. Sie können jeden Namen verwenden, der für Sie von Bedeutung ist. Die einzige Voraussetzung ist, dass die Zeichenfolge zwischen dem Task Runner und den entsprechenden Pipeline-Aktivitäten übereinstimmt. Der Name der Auftragnehmergruppe ist an eine Region gebunden. Auch wenn es in anderen Regionen identische Namen für Arbeitsgruppen gibt, ruft Task Runner immer Aufgaben aus der Region ab, die in angegeben ist--region.</p>
--taskrunnerId	Die ID des Task-Runners für die Berichterstellung zum Fortschritt. Optional.
--output	Das Task Runner-Verzeichnis für Protokollausgabedateien. Optional. Protokolldateien werden in einem lokalen Verzeichnis gespeichert, bis sie an Amazon S3 übertragen werden. Diese Option überschreibt das Standard-Verzeichnis.

Befehlszeilen-Parameter	Beschreibung
--region	<p>Die -Region, die verwendet werden soll. Optional, aber es wird empfohlen, die Region immer festzulegen. Wenn Sie die Region nicht angeben, ruft Task Runner Aufgaben aus der Standard-Serviceregion ab. <code>us-east-1</code></p> <p>Andere unterstützte Regionen sind: <code>eu-west-1</code> , <code>ap-northeast-1</code> , <code>ap-southeast-2</code> , <code>us-west-2</code> .</p>
--logUri	<p>Der Amazon S3 S3-Zielpfad für Task Runner, um Protokolldateien stündlich zu sichern. Wenn Task Runner beendet wird, werden aktive Protokolle im lokalen Verzeichnis in den Amazon S3 S3-Zielordner übertragen.</p>
--proxyHost	<p>Der Host des Proxys, der von Task Runner-Clients verwendet wird, um eine Verbindung zu AWS-Services herzustellen.</p>
--proxyPort	<p>Port des Proxy-Hosts, der von Task Runner-Clients verwendet wird, um eine Verbindung zu AWS-Services herzustellen.</p>
--proxyUsername	<p>Der Benutzername für den Proxy</p>
--proxyPassword	<p>Das Passwort für den Proxy.</p>
--proxyDomain	<p>Der Windows-Domänenname für NTLM Proxy.</p>
--proxyWorkstation	<p>Der Windows-Arbeitsstationsname für NTLM Proxy.</p>

Task-Runner mit einem Proxy verwenden

Wenn Sie einen Proxy-Host verwenden, können Sie beim Aufrufen von Task-Runner entweder seine [Konfiguration](#) angeben oder die Umgebungsvariable HTTPS_PROXY festlegen. Die mit Task-Runner verwendete Umgebungsvariable akzeptiert dieselbe Konfiguration, die für die [AWS-Befehlszeilenschnittstelle](#) verwendet wird.

Task Runner und Benutzerdefiniert AMIs

Wenn Sie ein Ec2Resource Objekt für Ihre Pipeline angeben, AWS Data Pipeline wird mithilfe eines AMI, das Task Runner für Sie installiert und konfiguriert, eine EC2 Instanz für Sie erstellt. In diesem Fall ist ein PV-kompatibler Instancetyp erforderlich. Alternativ können Sie mit Task Runner ein benutzerdefiniertes AMI erstellen und dann die ID dieses AMI mithilfe des imageId Felds des Ec2Resource Objekts angeben. Weitere Informationen finden Sie unter [Ec2Resource](#).

Ein benutzerdefiniertes AMI muss die folgenden Anforderungen erfüllen AWS Data Pipeline , um es erfolgreich für Task Runner verwenden zu können:

- Erstellen Sie das AMI in derselben Region, in der die Instances ausgeführt werden. Weitere Informationen finden Sie unter [Creating Your Own AMI](#) im EC2 Amazon-Benutzerhandbuch.
- Stellen Sie sicher, dass der Virtualisierungstyp des AMI vom Instancetyp unterstützt wird, den Sie verwenden möchten. Beispielsweise benötigen die Instancetypen I2 und G2 einen HVM AMI, und die Instancetypen T1, C1, M1 und M2 erfordern einen PV AMI. Weitere Informationen finden Sie unter [Linux AMI Virtualization Types](#) im EC2 Amazon-Benutzerhandbuch.
- Installieren Sie die folgenden Software:
 - Linux
 - Bash
 - wget
 - unzip
 - Java 1.6 oder 1.8
 - cloud-init
- Erstellen und konfigurieren Sie einen Benutzer mit dem Namenec2-user.

Fehlerbehebung

Wenn Sie ein Problem mit haben AWS Data Pipeline, ist das häufigste Symptom, dass eine Pipeline nicht läuft. Sie können die in der Konsole und der Befehlszeile angezeigten Informationen verwenden, um das Problem zu bestimmen und nach einer Lösung zu suchen.

Inhalt

- [Suchen von Fehlern in Pipelines](#)
- [Identifizieren des Amazon EMR-Clusters, der Ihre Pipeline bedient](#)
- [Interpretieren der Pipeline-Statusdetails](#)
- [Lokalisieren von Fehlerprotokollen](#)
- [Beheben typischer Probleme](#)

Suchen von Fehlern in Pipelines

Die AWS Data Pipeline Konsole ist ein praktisches Tool, mit dem Sie den Status Ihrer Pipelines visuell überwachen und Fehler im Zusammenhang mit fehlgeschlagenen oder unvollständigen Pipeline-Läufen leicht lokalisieren können.

So suchen Sie nach Fehlerinformationen zu fehlgeschlagenen oder unvollständigen Pipeline-Ausführungen.

1. Wenn auf der Seite List Pipelines in der Spalte Status einer Pipeline-Instance ein anderer Status als FINISHED angezeigt wird, wartet die betreffende Pipeline auf die Erfüllung einer Vorbedingung oder es ist ein Problem mit der Pipeline aufgetreten.
2. Suchen Sie auf der Seite List Pipelines (Pipelines auflisten) die Instance-Pipeline und klicken Sie auf das Dreieck links daneben, um den Detailbereich zu erweitern.
3. Klicken Sie unten in diesem Feld auf View execution details (Ausführungsdetails anzeigen); das Feld Instance summary (Instance-Zusammenfassung) wird geöffnet, um die Details der ausgewählten Instance anzuzeigen.
4. Klicken Sie im Bereich Instance summary (Instance-Zusammenfassung) auf das Dreieck neben der Instance, um zusätzliche Details zur Instance anzuzeigen, und wählen Sie Details, More.... Wenn der Status der ausgewählten Instance FAILED ist, enthält das Detailfeld Einträge für die Fehlermeldung, das errorStackTrace und weitere Informationen. Sie können diese Informationen in einer Datei speichern. Wählen Sie OK aus.

5. Klicken Sie im Bereich Instance summary (Instance-Zusammenfassung) auf Attempts (Versuche), um Details für jede Versuchszeile anzuzeigen.
6. Um eine Aktion mit der unvollständigen oder fehlgeschlagenen Instance durchzuführen, aktivieren Sie das Kontrollkästchen neben der Instance. Hiermit aktivieren Sie die Aktionen. Wählen Sie dann eine Aktion (Rerun | Cancel | Mark Finished).

Identifizieren des Amazon EMR-Clusters, der Ihre Pipeline bedient

Wenn ein EMRCluster oder EMRActivity fehlschlägt und die von der AWS Data Pipeline Konsole bereitgestellten Fehlerinformationen unklar sind, können Sie den Amazon EMR-Cluster, der Ihre Pipeline bedient, mithilfe der Amazon EMR-Konsole identifizieren. Auf diese Weise können Sie die von Amazon EMR bereitgestellten Protokolle finden, um weitere Informationen zu aufgetretenen Fehlern zu erhalten.

Um detailliertere Amazon EMR-Fehlerinformationen zu sehen

1. Wählen Sie in der AWS Data Pipeline Konsole das Dreieck neben der Pipeline-Instance aus, um die Instance-Details zu erweitern.
2. Klicken Sie auf View execution details (Ausführungsdetails anzeigen) und dann auf das Dreieck neben der Komponente.
3. Klicken Sie in der Spalte Details auf More... (Mehr...). Der Informationsbildschirm wird geöffnet und zeigt eine Liste der Details der Komponente an. Suchen und kopieren Sie den instanceParent-Wert auf dem Bildschirm, z. B.:
@EmrActivityId_xiFDD_2017-09-30T21:40:13
4. Navigieren Sie zur Amazon EMR-Konsole, suchen Sie nach einem Cluster mit dem entsprechenden InstanceParent-Wert im Namen und wählen Sie dann Debug aus.

 Note

Damit die Debug-Schaltfläche funktioniert, muss Ihre Pipeline-Definition die Option auf true und die EmrActivity enableDebugging EmrLogUri Option auf einen gültigen Pfad gesetzt haben.

5. Nachdem Sie nun wissen, welcher Amazon EMR-Cluster den Fehler enthält, der Ihren Pipeline-Ausfall verursacht hat, folgen Sie den [Tipps zur Fehlerbehebung](#) im Amazon EMR Developer Guide.

Interpretieren der Pipeline-Statusdetails

Die verschiedenen Statusstufen, die in der AWS Data Pipeline Konsole und der CLI angezeigt werden, geben den Zustand einer Pipeline und ihrer Komponenten an. Der Pipeline-Status ist vereinfacht ausgedrückt ein Überblick über eine Pipeline. Wenn Sie weitere Informationen benötigen, zeigen Sie den Status der einzelnen Pipeline-Komponenten an. Sie können dazu in der Konsole auf die Pipeline und deren Komponenten klicken oder die Details der Pipeline-Komponenten über die CLI abrufen.

Statuscodes

ACTIVATING

Die Komponente oder Ressource wird gestartet, z. B. eine EC2 Instanz.

CANCELED

Die Komponente wurde von einem Benutzer abgebrochen oder AWS Data Pipeline bevor sie ausgeführt werden konnte. Dies kann automatisch geschehen, wenn ein Fehler in einer anderen Komponente oder Ressource auftritt, von der diese Komponente abhängt.

CASCADE_FAILED

Die Komponente oder Ressource wurde aufgrund eines Kaskadenfehlers aufgrund einer ihrer Abhängigkeiten storniert, aber die Komponente war wahrscheinlich nicht die ursprüngliche Ursache des Fehlers.

DEACTIVATING

Die Pipeline wird deaktiviert.

FAILED

Bei der Komponente oder Ressource ist ein Fehler aufgetreten und sie funktioniert nicht mehr. Wenn eine Komponente oder Ressource ausfällt, kann dies dazu führen, dass Stornierungen und Fehler auf Komponenten übertragen werden, die von ihr abhängig sind.

FINISHED

Die Komponente hat die ihr zugewiesene Arbeit abgeschlossen.

INACTIVE

Die Pipeline wurde deaktiviert.

PAUSED

Die Komponente wurde angehalten und führt derzeit ihre Arbeit nicht aus.

PENDING

Die Pipeline ist bereit, zum ersten Mal aktiviert zu werden.

RUNNING

Die Ressource läuft und ist bereit, Arbeit zu empfangen.

SCHEDULED

Die Ausführung der Ressource ist geplant.

SHUTTING_DOWN

Die Ressource wird nach erfolgreichem Abschluss ihrer Arbeit heruntergefahren.

SKIPPED

Die Komponente hat Ausführungsintervalle nach der Aktivierung der Pipeline übersprungen, wobei ein Zeitstempel verwendet wurde, der nach dem aktuellen Zeitplan liegt.

TIMEDOUT

Die Ressource hat den `terminateAfter` Schwellenwert überschritten und wurde von AWS Data Pipeline gestoppt. Wenn die Ressource diesen Status erreicht hat, werden die `retryTimeout` Werte `actionOnResourceFailureRetryDelay`, und für diese Ressource AWS Data Pipeline ignoriert. Dieser Status gilt nur für Ressourcen.

VALIDATING

Die Pipeline-Definition wird von validiert AWS Data Pipeline.

WAITING_FOR_RUNNER

Die Komponente wartet darauf, dass ihr Worker-Client ein Arbeitselement abruft. Die Beziehung zwischen Komponente und Worker-Client wird durch die `workerGroup` Felder `runsOn` oder gesteuert, die von dieser Komponente definiert sind.

WAITING_ON_DEPENDENCIES

Die Komponente überprüft, ob ihre standardmäßigen und vom Benutzer konfigurierten Voraussetzungen erfüllt sind, bevor sie ihre Arbeit ausführt.

Lokalisieren von Fehlerprotokollen

In diesem Abschnitt wird erklärt, wie Sie die verschiedenen Protokolle finden, die AWS Data Pipeline geschrieben werden. Anhand dieser Informationen können Sie die Ursache bestimmter Ausfälle und Fehler ermitteln.

Pipeline-Protokolle

Wir empfehlen, Pipelines so zu konfigurieren, dass Protokolldateien an einem persistenten Speicherort erstellt werden, wie im folgenden Beispiel, in dem Sie das `pipelineLogUri` Feld für das `Default` Objekt einer Pipeline verwenden, um zu veranlassen, dass alle Pipeline-Komponenten standardmäßig einen Amazon S3 S3-Protokollspeicherort verwenden (Sie können dies überschreiben, indem Sie einen Protokollspeicherort in einer bestimmten Pipeline-Komponente konfigurieren).

Note

Task Runner speichert seine Protokolle standardmäßig an einem anderen Speicherort, der möglicherweise nicht verfügbar ist, wenn die Pipeline abgeschlossen ist und die Instanz, auf der Task Runner ausgeführt wird, beendet wird. Weitere Informationen finden Sie unter [Die Task Runner-Protokollierung wird überprüft](#).

Um den Protokollspeicherort mithilfe der AWS Data Pipeline CLI in einer Pipeline-JSON-Datei zu konfigurieren, beginnen Sie Ihre Pipeline-Datei mit dem folgenden Text:

```
{ "objects": [ { "id": "Default", "pipelineLogUri": "s3://amzn-s3-demo-bucket/error_logs" }, ... ] }
```

Nachdem Sie ein Pipeline-Protokollverzeichnis konfiguriert haben, erstellt Task Runner eine Kopie der Protokolle in Ihrem Verzeichnis mit derselben Formatierung und denselben Dateinamen, die im vorherigen Abschnitt über Task Runner-Protokolle beschrieben wurden.

Hadoop-Job- und Amazon EMR-Schrittprotokolle

Bei jeder Hadoop-basierten Aktivität wie [HadoopActivity](#), [HiveActivity](#), oder können [PigActivity](#) Sie Hadoop-Jobprotokolle an dem Ort einsehen, der im Runtime-Slot, zurückgegeben wurde. hadoopJobLog [EmrActivity](#) verfügt über eigene Protokollierungsfunktionen, und diese Protokolle werden an dem von Amazon EMR ausgewählten Speicherort gespeichert und vom Runtime-Slot zurückgegeben. emrStepLog Weitere Informationen finden Sie unter [Anzeigen von Protokolldateien](#) im Amazon EMR-Entwicklerhandbuch.

Beheben typischer Probleme

In diesem Thema werden verschiedene Symptome von AWS Data Pipeline Problemen und die empfohlenen Schritte zu ihrer Lösung beschrieben.

Inhalt

- [Pipeline bleibt im Status PENDING](#)
- [Pipeline-Komponente bleibt im Status WAITING_FOR_RUNNER](#)
- [Pipeline-Komponente bleibt im Status WAITING_ON_DEPENDENCIES](#)
- [Ausführung beginnt nicht zum geplanten Zeitpunkt](#)
- [Pipeline-Komponenten werden in der falschen Reihenfolge ausgeführt](#)
- [EMR-Cluster schlägt mit Fehlermeldung fehl: The security token included in the request is invalid](#)
- [Unzureichende Berechtigungen für den Zugriff auf Ressourcen](#)
- [Statuscode: 400 Fehlercode: PipelineNotFoundException](#)
- [Pipeline-Erstellung führt zu einem Sicherheits-Token-Fehler](#)
- [Pipeline-Details werden nicht in der Konsole angezeigt](#)
- [Error in remote runner Status Code: 404, AWS Service: Amazon S3](#)
- [Access Denied - Not Authorized to Perform Function datapipeline:](#)
- [Ältere Versionen von Amazon EMR AMIs können falsche Daten für große CSV-Dateien erstellen](#)
- [AWS Data Pipeline Steigende Grenzwerte](#)

Pipeline bleibt im Status PENDING

Wenn eine Pipeline dauerhaft im Status PENDING bleibt, weist dies darauf hin, dass sie noch nicht aktiviert wurde oder die Aktivierung aufgrund eines Fehlers in der Pipeline-Definition fehlgeschlagen

ist. Stellen Sie sicher, dass Sie keine Fehler erhalten haben, als Sie Ihre Pipeline über die AWS Data Pipeline CLI eingereicht haben oder als Sie versucht haben, Ihre Pipeline über die AWS Data Pipeline Konsole zu speichern oder zu aktivieren. Überprüfen Sie außerdem, ob die Definition der Pipeline gültig ist.

So zeigen Sie die Pipeline-Definition unter Verwendung der Befehlszeile an:

```
aws datapipeline --get-pipeline-definition --pipeline-id df-EXAMPLE_PIPELINE_ID
```

Überzeugen Sie sich, dass die Definition vollständig ist und keine Syntaxfehler enthält. Achten Sie insbesondere darauf, dass keine schließenden Klammern, erforderlichen Kommas oder Verweise fehlen. Wir empfehlen, einen Texteditor zu verwenden, der die Syntax von JSON-Dateien visuell darstellen und validieren kann.

Pipeline-Komponente bleibt im Status WAITING_FOR_RUNNER

Wenn sich die Pipeline im Status SCHEDULED befindet und einzelne oder mehrere Aufgaben dauerhaft im Status WAITING_FOR_RUNNER bleiben, stellen Sie sicher, dass das Feld "runsOn" oder "workerGroup" dieser Aufgaben einen gültigen Wert enthält. Falls beide Werte leer sind oder fehlen, kann die betreffende Aufgabe nicht gestartet werden, da es keine Zuordnung zwischen ihr und dem Worker zur Durchführung der Aufgaben gibt. In diese Situation haben Sie zwar durchzuführende Aufgaben definiert, aber nicht festgelegt, auf welchem Computer dies geschehen soll. Stellen Sie gegebenenfalls sicher, dass der der Pipeline-Komponente zugewiesene WorkerGroup-Wert genau den gleichen Namen und die gleiche Groß- und Kleinschreibung hat wie der WorkerGroup-Wert, den Sie für Task Runner konfiguriert haben.

Note

Wenn Sie einen runsOn-Wert angeben und workerGroup vorhanden ist, wird workerGroup ignoriert.

Eine weitere mögliche Ursache für dieses Problem ist, dass der Endpunkt und der Zugriffsschlüssel, die Task Runner zur Verfügung gestellt werden, nicht mit der AWS Data Pipeline Konsole oder dem Computer identisch sind, auf dem die AWS Data Pipeline CLI-Tools installiert sind. Möglicherweise haben Sie neue Pipelines ohne sichtbare Fehler erstellt, aber Task Runner fragt aufgrund der unterschiedlichen Anmeldeinformationen den falschen Standort ab oder fragt den richtigen

Standort ab, der nicht über ausreichende Berechtigungen verfügt, um die in der Pipeline-Definition angegebene Arbeit zu identifizieren und auszuführen.

Pipeline-Komponente bleibt im Status WAITING_ON_DEPENDENCIES

Wenn sich die Pipeline im Status SCHEDULED befindet und einzelne oder mehrere Aufgaben dauerhaft im Status WAITING_ON_DEPENDENCIES bleiben, vergewissern Sie sich, dass die Vorbedingungen erfüllt sind. Wenn die Vorbedingungen des ersten Objekts in der Logikkette nicht erfüllt sind, kann keines der Objekte, die vom ersten Objekt abhängig sind, den Status WAITING_ON_DEPENDENCIES verlassen.

Sehen Sie sich als Beispiel den folgenden Auszug aus einer Pipeline-Definition an. In diesem Fall hat das InputData Objekt die Vorbedingung „Bereit“, die angibt, dass die Daten vorhanden sein müssen, bevor das InputData Objekt fertiggestellt ist. Wenn die Daten nicht existieren, verbleibt das InputData Objekt im WAITING_ON_DEPENDENCIES Status und wartet darauf, dass die im Pfadfeld angegebenen Daten verfügbar werden. Alle Objekte, die von Elke InputData abhängig sind, verbleiben in einem WAITING_ON_DEPENDENCIES Zustand und warten darauf, dass das InputData Objekt den FINISHED Status erreicht.

```
{  
    "id": "InputData",  
    "type": "S3DataNode",  
    "filePath": "s3://elasticmapreduce/samples/wordcount/wordSplitter.py",  
    "schedule": {"ref": "MySchedule"},  
    "precondition": "Ready"  
},  
{  
    "id": "Ready",  
    "type": "Exists"  
...  
}
```

Überprüfen Sie außerdem, ob die Objekte über die Berechtigungen zum Zugriff auf die Daten verfügen. Wenn im vorherigen Beispiel die Informationen im Feld mit den Anmeldeinformationen nicht berechtigt waren, auf die im Pfadfeld angegebenen Daten zuzugreifen, würde das InputData Objekt in einem WAITING_ON_DEPENDENCIES Status stecken bleiben, weil es nicht auf die im Pfadfeld angegebenen Daten zugreifen kann, selbst wenn diese Daten existieren.

Es ist auch möglich, dass einer Ressource, die mit Amazon S3 kommuniziert, keine öffentliche IP-Adresse zugeordnet ist. So muss beispielsweise eine Ec2Resource in einem öffentlichen Subnetz über eine öffentliche IP-Adresse verfügen.

Und schließlich können Ressourceninstanzen unter bestimmten Bedingungen den Status WAITING_ON_DEPENDENCIES viel früher erreichen als ihre zugeordneten Aktivitäten, die zur Ausführung geplant sind. Dies kann den Eindruck erwecken, dass die Ressource oder Aktivität fehlschlägt.

Ausführung beginnt nicht zum geplanten Zeitpunkt

Vergewissern Sie sich, dass Sie den richtigen Zeitplantyp ausgewählt haben, der festlegt, ob die Aufgabe am Anfang des Zeitplanintervalls (Cron-Zeitplantyp) oder an dessen Ende (Zeitreihen-Zeitplantyp) ausgeführt wird.

Vergewissern Sie sich außerdem, dass Sie die Datumsangaben in Ihren Zeitplanobjekten korrekt angegeben haben startDateTime und dass die endDateTime Werte und im UTC-Format vorliegen, wie im folgenden Beispiel:

```
{  
    "id": "MySchedule",  
    "startDateTime": "2012-11-12T19:30:00",  
    "endDateTime": "2012-11-12T20:30:00",  
    "period": "1 Hour",  
    "type": "Schedule"  
},
```

Pipeline-Komponenten werden in der falschen Reihenfolge ausgeführt

Möglicherweise stellen Sie fest, dass die Pipeline-Komponenten nicht in der durch die Start- und Endzeiten festgelegten Reihenfolge oder nicht in der erwarteten Abfolge ausgeführt werden. Sie müssen wissen, dass Pipeline-Komponenten gleichzeitig gestartet werden können, wenn ihre Vorbedingungen zum Ausführungszeitpunkt erfüllt sind. Anders ausgedrückt: Die Pipeline-Komponenten werden nicht standardmäßig der Reihe nach ausgeführt. Wenn Sie eine bestimmte Ausführungsreihenfolge verwenden möchten, müssen Sie diese mithilfe von Vorbedingungen und dependsOn-Feldern festlegen.

Stellen Sie außerdem sicher, dass das dependsOn-Feld einen Verweis auf die richtigen Vorbedingungs-Pipeline-Komponenten enthält und dass alle erforderlichen Zeiger zwischen den Komponenten vorhanden sind, um die gewünschte Reihenfolge zu erreichen.

EMR-Cluster schlägt mit Fehlermeldung fehl: The security token included in the request is invalid

Überprüfen Sie Ihre IAM-Rollen, Richtlinien und Vertrauensbeziehungen wie unter beschrieben [IAM-Rollen für AWS Data Pipeline](#).

Unzureichende Berechtigungen für den Zugriff auf Ressourcen

Berechtigungen, die Sie für IAM-Rollen festlegen, bestimmen, ob Sie auf Ihre EMR-Cluster und EC2 -Instances zugreifen AWS Data Pipeline können, um Ihre Pipelines auszuführen. Darüber hinaus bietet IAM das Konzept von Vertrauensbeziehungen, das noch weiter geht und die Erstellung von Ressourcen in Ihrem Namen ermöglicht. Wenn Sie beispielsweise eine Pipeline erstellen, die eine EC2 Instanz verwendet, um einen Befehl zum Verschieben von Daten auszuführen, AWS Data Pipeline kann ich diese EC2 Instanz für Sie bereitstellen. Wenn Sie auf Probleme stoßen, insbesondere wenn es um Ressourcen geht, auf die Sie manuell zugreifen können, auf die Sie jedoch AWS Data Pipeline nicht zugreifen können, überprüfen Sie Ihre IAM-Rollen, Richtlinien und Vertrauensbeziehungen wie unter beschrieben. [IAM-Rollen für AWS Data Pipeline](#)

Statuscode: 400 Fehlercode: PipelineNotFoundException

Dieser Fehler bedeutet, dass Ihre IAM-Standardrollen möglicherweise nicht über die erforderlichen Berechtigungen verfügen, um ordnungsgemäß AWS Data Pipeline zu funktionieren. Weitere Informationen finden Sie unter [IAM-Rollen für AWS Data Pipeline](#).

Pipeline-Erstellung führt zu einem Sicherheits-Token-Fehler

Wenn Sie versuchen, eine Pipeline zu erstellen, wird die folgende Fehlermeldung angezeigt:

Failed to create pipeline with 'pipeline_name'. Fehler: UnrecognizedClientException Das Sicherheits-Token der Anfrage ist ungültig.

Pipeline-Details werden nicht in der Konsole angezeigt

Der AWS Data Pipeline Konsolen-Pipeline-Filter gilt für das geplante Startdatum für eine Pipeline, unabhängig davon, wann die Pipeline eingereicht wurde. Es ist möglich, eine neue Pipeline mit einem geplanten Anfangsdatum zu übermitteln, das in der Vergangenheit liegt. Diese Pipeline wird dann vom Standardfilter gefiltert und nicht angezeigt. Um die Details der Pipeline anzuzeigen, ändern Sie den Datumsfilter so, dass das geplante Startdatum der Pipeline innerhalb des Datumsbereichs des Filters liegt.

Error in remote runner Status Code: 404, AWS Service: Amazon S3

Dieser Fehler bedeutet, dass Task Runner nicht auf Ihre Dateien in Amazon S3 zugreifen konnte. Vergewissern Sie sich, dass folgende Bedingungen erfüllt sind:

- Die Anmeldeinformationen wurden richtig festgelegt.
- Der Amazon S3 S3-Bucket, auf den Sie zugreifen möchten, ist vorhanden
- Sie sind für den Zugriff auf den Amazon S3 S3-Bucket autorisiert

Access Denied - Not Authorized to Perform Function datapipeline:

In den Task Runner-Protokollen wird möglicherweise ein Fehler angezeigt, der dem folgenden ähnelt:

- ERROR Status Code: 403
- AWS-Dienst: DataPipeline
- AWS-Fehlercode: AccessDenied
- AWS-Fehlermeldung: Benutzer: arn:aws:sts: :XXXXXXXXXXXXX:Federated-User/i-XXXXXXX ist nicht zur Durchführung autorisiert: datapipeline: PollForTask

 Note

Kann PollForTask in dieser Fehlermeldung durch Namen anderer Berechtigungen ersetzt AWS Data Pipeline werden.

Diese Fehlermeldung weist darauf hin, dass die von Ihnen angegebene IAM-Rolle zusätzliche Berechtigungen benötigt, um mit AWS Data Pipeline ihr interagieren zu können. Stellen Sie sicher, dass Ihre IAM-Rollenrichtlinie die folgenden Zeilen enthält, wobei PollForTask diese durch den Namen der Berechtigung ersetzt werden, die Sie hinzufügen möchten (verwenden Sie *, um alle Berechtigungen zu gewähren). Weitere Informationen zum Erstellen einer neuen IAM-Rolle und zum Anwenden einer Richtlinie auf diese Rolle finden Sie unter [Verwaltung von IAM-Richtlinien](#) im Handbuch Using IAM.

```
{  
  "Action": [ "datapipeline:PollForTask" ],  
  "Effect": "Allow",
```

```
"Resource": ["*"]  
}
```

Ältere Versionen von Amazon EMR AMIs können falsche Daten für große CSV-Dateien erstellen

Auf Amazon EMR vor 3.9 (3.8 und darunter) AMIs wurde ein benutzerdefinierter Code InputFormat zum Lesen und Schreiben von CSV-Dateien zur Verwendung mit MapReduce Jobs AWS Data Pipeline verwendet. Dies wird verwendet, wenn der Service Tabellen zu und von Amazon S3 bereitstellt. Dabei InputFormat wurde ein Problem entdeckt, bei dem das Lesen von Datensätzen aus großen CSV-Dateien dazu führen kann, dass Tabellen erstellt werden, die nicht korrekt kopiert werden. Dieses Problem wurde in späteren Amazon EMR-Versionen behoben. Verwenden Sie Amazon EMR-AMI 3.9 oder eine Amazon EMR-Version 4.0.0 oder höher.

AWS Data Pipeline Steigende Grenzwerte

Gelegentlich können Sie bestimmte AWS Data Pipeline Systemgrenzen überschreiten. So können beispielsweise höchstens 20 Pipelines mit jeweils 50 Objekten erstellt werden. Falls Sie mehr Pipelines benötigen, können Sie mehrere Pipelines zusammenführen. Sie erhalten dann weniger Pipelines mit mehr Objekten in jeder. Weitere Informationen zu den Limits für AWS Data Pipeline finden Sie unter [AWS Data Pipeline Grenzwerte](#). Sollte diese Abhilfemaßnahme nicht zum Erfolg führen, können Sie mit dem folgenden Formular eine Kapazitätserhöhung anfordern: [Erhöhen des Pipeline-Limits](#).

AWS Data Pipeline Grenzwerte

AWS Data Pipeline schränkt die Ressourcen, die Sie zuweisen können, und die Geschwindigkeit, mit der Sie Ressourcen zuweisen können, ein, um sicherzustellen, dass für alle Benutzer genügend Kapazität zur Verfügung steht.

Inhalt

- [Kontolimits](#)
- [Limits für Webservice-Aufrufe](#)
- [Überlegungen zur Skalierung](#)

Kontolimits

Die folgenden Beschränkungen gelten für ein einzelnes Konto. AWS Wenn Sie zusätzliche Kapazität benötigen, können Sie das [Antragsformular für das Amazon Web Services Support Center](#) verwenden, um Ihre Kapazität zu erhöhen.

Attribut	Limit	Einstellbar
Anzahl Pipelines	100	Ja
Anzahl Objekte pro Pipeline	100	Ja
Anzahl aktiver Instances pro Objekt	5	Ja
Anzahl Felder pro Objekt	50	Nein
Anzahl der UTF8 Byte pro Feldname oder ID	256	Nein
Anzahl der UTF8 Byte pro Feld	10,240	Nein

Attribut	Limit	Einstellbar
Anzahl der UTF8 Byte pro Objekt	15.360 (einschl. Feldnamen)	Nein
Erstellungsrate einer Instance von einem Objekt	1 pro 5 Minuten	Nein
Neuversuche einer Pipeline-Aktivität	5 pro Aufgabe	Nein
Minimale Verzögerung zwischen Neuversuchen	2 Minuten	Nein
Minimales Planungintervall	15 Minuten	Nein
Maximale Anzahl Aggregationen zu einem Objekt	32	Nein
Maximale Anzahl von EC2 Instanzen pro Ec2Resource-Objekt	1	Nein

Limits für Webservice-Aufrufe

AWS Data Pipeline begrenzt die Geschwindigkeit, mit der Sie die Webdienst-API aufrufen können. Diese Beschränkungen gelten auch für AWS Data Pipeline Agenten, die die Webservice-API in Ihrem Namen aufrufen, z. B. die Konsole, CLI und Task Runner.

Die folgenden Beschränkungen gelten für ein einzelnes AWS Konto. Die Gesamtnutzung des Kontos, einschließlich der Nutzung durch -Benutzer, kann diese Grenzwerte also nicht überschreiten.

Mit der Burst-Rate können Sie Webservice-Aufrufe in inaktiven Zeiträumen einsparen und sie alle in einem kurzen Zeitraum aufbrauchen. CreatePipeline Hat beispielsweise eine reguläre Rate von

einem Anruf alle fünf Sekunden. Wenn Sie den Service 30 Sekunden nicht aufrufen, haben Sie 6 Aufrufe gespart. Sie können dann den Webservice sechsmal in einer Sekunde aufrufen. Da dieser Wert unter dem Burst-Limit liegt und Ihre durchschnittlichen Aufrufe auf dem regulären Ratenlimit belässt, werden die Aufrufe nicht gedrosselt.

Wenn Sie das Raten- und das Burst-Limit überschreiten, schlägt der Webservice-Aufruf fehl und gibt eine Drosselungsausnahme zurück. Die Standardimplementierung eines Workers, Task Runner, wiederholt API-Aufrufe, die fehlschlagen, automatisch mit einer Drosselungsausnahme. Task Runner hat einen Back-off, sodass nachfolgende Versuche, die API aufzurufen, in immer längeren Intervallen erfolgen. Wenn Sie einen Worker schreiben, empfehlen wir, dass Sie eine ähnliche Logik für wiederholte Versuche implementieren.

Diese Limits gelten für ein einzelnes AWS Konto.

API	Reguläres Ratenlimit	Burst-Limit
ActivatePipeline	1 Aufruf pro Sekunde	100 Aufrufe
CreatePipeline	1 Aufruf pro Sekunde	100 Aufrufe
DeletePipeline	1 Aufruf pro Sekunde	100 Aufrufe
DescribeObjects	2 Aufrufe pro Sekunde	100 Aufrufe
DescribePipelines	1 Aufruf pro Sekunde	100 Aufrufe
GetPipelineDefinition	1 Aufruf pro Sekunde	100 Aufrufe
PollForTask	2 Aufrufe pro Sekunde	100 Aufrufe
ListPipelines	1 Aufruf pro Sekunde	100 Aufrufe
PutPipelineDefinition	1 Aufruf pro Sekunde	100 Aufrufe
QueryObjects	2 Aufrufe pro Sekunde	100 Aufrufe
ReportTaskProgress	10 Aufrufe pro Sekunde	100 Aufrufe
SetTaskStatus	10 Aufrufe pro Sekunde	100 Aufrufe
SetStatus	1 Aufruf pro Sekunde	100 Aufrufe

API	Reguläres Ratenlimit	Burst-Limit
ReportTaskRunnerHeartbeat	1 Aufruf pro Sekunde	100 Aufrufe
ValidatePipelineDefinition	1 Aufruf pro Sekunde	100 Aufrufe

Überlegungen zur Skalierung

AWS Data Pipeline skaliert, um eine große Anzahl gleichzeitiger Aufgaben zu bewältigen, und Sie können es so konfigurieren, dass automatisch die Ressourcen bereitgestellt werden, die für die Bewältigung großer Workloads erforderlich sind. Diese automatisch erstellten Ressourcen sind von Ihnen steuerbar und werden für die Ressourcenlimits für Ihr AWS -Konto berücksichtigt. Wenn Sie beispielsweise so konfigurieren, AWS Data Pipeline dass automatisch ein Amazon EMR-Cluster mit 20 Knoten zur Verarbeitung von Daten erstellt wird und für Ihr AWS Konto ein EC2 Instance-Limit von 20 festgelegt ist, können Sie versehentlich Ihre verfügbaren Backfill-Ressourcen erschöpfen. Daher sollten Sie diese Ressourceneinschränkungen bei Ihrem Design berücksichtigen oder Ihre Kontolimits entsprechend erweitern.

Wenn Sie zusätzliche Kapazität benötigen, können Sie das [Antragsformular für das Amazon Web Services Support Center](#) verwenden, um Ihre Kapazität zu erhöhen.

AWS Data Pipeline Ressourcen

Im Folgenden finden Sie Ressourcen für die Verwendung von AWS Data Pipeline.

- [AWS Data Pipeline Produktinformationen](#) — Die primäre Webseite mit Informationen über AWS Data Pipeline.
- [AWS Data Pipeline Häufig gestellte technische Fragen](#) — Enthält die 20 häufigsten Fragen, die Entwickler zu diesem Produkt stellen.
- [Versionshinweise](#) — Bieten einen allgemeinen Überblick über die aktuelle Version. Im Einzelnen werden neue Funktionen, Korrekturen und bekannte Probleme vorgestellt.
- [Diskussionsforen für AWS Data Pipeline](#) — Ein Community-Forum für Entwickler, um technische Fragen zu Amazon Web Services zu erörtern.
- [Kurse und Workshops — Links zu rollen-](#) und Spezialkursen sowie zu Übungen zum Selbststudium, mit denen Sie Ihre Fähigkeiten verbessern und praktische Erfahrungen sammeln können. AWS
- [AWS Developer Center](#) — Erkunden Sie Tutorials, laden Sie Tools herunter und erfahren Sie mehr über Veranstaltungen für Entwickler. AWS
- [AWS Entwickertools](#) — Links zu Entwickertools SDKs, IDE-Toolkits und Befehlszeilentools für die Entwicklung und Verwaltung von AWS Anwendungen.
- [Ressourcencenter für die ersten Schritte](#) — Erfahren Sie AWS-Konto, wie Sie Ihre erste Anwendung einrichten, der AWS Community beitreten und sie starten.
- [Praktische Tutorials](#) — Folgen Sie den step-by-step Tutorials, um Ihre erste Anwendung zu starten. AWS
- [AWS Whitepapers](#) — Links zu einer umfassenden Liste von technischen AWS Whitepapers zu Themen wie Architektur, Sicherheit und Wirtschaft, die von Solutions Architects oder anderen technischen Experten verfasst wurden. AWS
- [AWS -Support Center](#) — Die zentrale Anlaufstelle für die Erstellung und Verwaltung Ihrer Fälle. AWS -Support Enthält auch Links zu anderen hilfreichen Ressourcen wie Foren, technischen FAQs Informationen, Servicestatus und AWS Trusted Advisor.
- [Support](#) — Die wichtigste Webseite mit Informationen über Support einen Support-Kanal mit schnellen Reaktionszeiten one-on-one, der Sie bei der Entwicklung und Ausführung von Anwendungen in der Cloud unterstützt.

- [Kontakt](#) – Zentraler Kontaktpunkt für Fragen zu AWS -Abrechnung, Konten, Ereignissen Missbrauch und anderen Problemen.
- [AWS Nutzungsbedingungen der Website](#) — Detaillierte Informationen zu unseren Urheberrechten und Marken, zu Ihrem Konto, Ihrer Lizenz und Ihrem Zugriff auf die Website sowie zu anderen Themen.

Dokumentverlauf

Diese Dokumentation ist mit der Version 2012-10-29 von verknüpft. AWS Data Pipeline

Änderung	Beschreibung	Veröffentlichungsdatum
AWS Data Pipeline ist für Neukunden nicht mehr verfügbar	AWS Data Pipeline ist für Neukunden nicht mehr verfügbar. Bestandskunden von AWS Data Pipeline können den Service weiterhin wie gewohnt nutzen. Weitere Informationen	25. Juli 2025
Dokumentation für die Durchführung bestimmter Verfahren mit der AWS CLI hinzugefügt. AWS Data Pipeline Prozeduren im Zusammenhang mit der Konsole wurden entfernt.	Weitere Informationen finden Sie unter Klonen Ihrer Pipeline , Anzeigen von Pipeline-Protokollen und Erstellen Sie mit der CLI eine Pipeline aus Data Pipeline-Vorlagen .	26. Mai 2023
Es wurden weitere Inhalte und Beispiele für die Migration von AWS Data Pipeline anderen alternativen Diensten hinzugefügt.	Das Thema für die Migration AWS Data Pipeline zu AWS Step Functions oder Amazon MWAA wurde mit weiteren Informationen zu den einzelnen Alternativen, Konzeptzuordnungen zwischen den Services und Beispielen aktualisiert. AWS Glue Weitere Informationen finden Sie unter Migration von Workloads von AWS Data Pipeline .	31. März 2023
Es wurden Informationen zur AWS Data Pipeline Unterstützung von hinzugefügtem IMDSv2.	AWS Data Pipeline unterstützt IMDSv2 Amazon EMR und EC2 Amazon-Ressourcen. Weitere Informationen finden Sie unter Datenschutz in AWS Data Pipeline , EmrCluster und Ec2Resource .	16. Dezember 2022

Änderung	Beschreibung	Veröffentlichungsdatum
Es wurde ein Thema für die Migration von AWS Data Pipeline anderen alternativen Diensten hinzugefügt.	Inzwischen gibt es andere AWS Dienste, die Kunden ein besseres Datenintegrationserlebnis bieten. Sie können typische Anwendungsfälle von entweder AWS Data Pipeline AWS Glue zu AWS Step Functions oder Amazon MWAA migrieren. Weitere Informationen finden Sie unter Migration von Workloads von AWS Data Pipeline .	16. Dezember 2022
Die Listen der unterstützten Amazon EC2 - und Amazon EMR-Instances wurden aktualisiert. Die Liste IDs der für die Instances AMIs verwendeten HVM (Hardware Virtual Machine) wurde aktualisiert.	Die Listen der unterstützten Amazon EC2 - und Amazon EMR-Instances wurden aktualisiert. Weitere Informationen finden Sie unter Unterstützte Instance-Typen für Pipeline-Aktivitäten . Die Liste IDs der für die Instances AMIs verwendeten HVM (Hardware Virtual Machine) wurde aktualisiert. Weitere Informationen finden Sie unter Syntax und bei der Suche nach <code>imageId</code> .	9. November 2018

Änderung	Beschreibung	Veröffentlichungsdatum
<p>Konfiguration für das Anhängen von Amazon EBS-Volumen an Clusterknoten und für den Start eines Amazon EMR-Clusters in einem privaten Subnetz hinzugefügt.</p>	<p>Konfigurationsoptionen wurden zu einem EMRcluster -Objekt hinzugefügt. Sie können diese Optionen in Pipelines verwenden, die Amazon EMR-Cluster verwenden.</p> <p>Verwenden Sie die TaskEbsConfiguration Felder <code>coreEbsConfiguration</code> <code>masterEbsConfiguration</code>, und, um das Anhängen von Amazon EBS-Volumes an Core-, Master- und Task-Knoten im Amazon EMR-Cluster zu konfigurieren. Weitere Informationen finden Sie unter EBS-Volumes zu Cluster-Knoten hinzufügen.</p> <p>Verwenden Sie die ServiceAccessSecurityGroupId Felder <code>emrManagedMasterSecurityGroupId</code> <code>emrManagedSlaveSecurityGroupId</code>, und, um einen Amazon EMR-Cluster in einem privaten Subnetz zu konfigurieren. Weitere Informationen finden Sie unter Einen Amazon EMR-Cluster in einem privaten Subnetz konfigurieren.</p> <p>Weitere Informationen zur EMRcluster -Syntax finden Sie unter EmrCluster.</p>	19. April 2018
<p>Die Liste der unterstützten Amazon EC2 - und Amazon EMR-Instances wurde hinzugefügt.</p>	<p>Die Liste der Instances wurde hinzugefügt, die standardmäßig AWS Data Pipeline erstellt wird, wenn Sie keinen Instance-Typ in der Pipeline-Definition angeben. Eine Liste der unterstützten Amazon EC2 - und Amazon EMR-Instances wurde hinzugefügt. Weitere Informationen finden Sie unter Unterstützte Instance-Typen für Pipeline-Aktivitäten.</p>	22. März 2018

Änderung	Beschreibung	Veröffentlichungsdatum
Unterstützung für On-Demand-Pipeline s wurde hinzugefügt.	<ul style="list-style-type: none"> Unterstützung für On-Demand-Pipelines hinzugefügt, die das Wiederholen einer Pipeline durch erneutes Aktivieren ermöglicht. 	22. Februar 2016
Zusätzliche Unterstützung für RDS-Datenbanken hinzugefügt	<ul style="list-style-type: none"> <code>rdsInstanceId</code>, <code>region</code> und <code>jdbcDrive rJarUri</code> wurden RdsDatabase hinzugefügt. <code>database</code> in SqlActivity wurde aktualisiert, um auch <code>RdsDatabase</code> zu unterstützen. 	17. August 2015
Zusätzliche JDBC-Unterstützung	<ul style="list-style-type: none"> <code>database</code> in SqlActivity wurde aktualisiert, um auch <code>JdbcDatabase</code> zu unterstützen. <code>jdbcDriverJarUri</code> wurde JdbcDatabase hinzugefügt. <code>initTimeout</code> wurde Ec2Resource und EmrCluste r hinzugefügt. Hinzufügung von <code>runAsUser</code> zu Ec2Resource. 	7. Juli 2015
HadoopActivity, Availability Zone und Spot-Support	<ul style="list-style-type: none"> Unterstützung der Übermittlung paralleler Arbeitsaufträge an Hadoop-Cluster wurde hinzugefügt. Weitere Informationen finden Sie unter HadoopActivity. Möglichkeit zur Anforderung von Spot-Instances mit Ec2Resource und EmrCluster hinzugefügt. Möglichkeit zum Starten von <code>EmrCluster</code> - Ressourcen in einer angegebenen Availability Zone wurden hinzugefügt. 	1. Juni 2015
Deaktivieren von Pipelines	Unterstützung zum Deaktivieren aktiver Pipelines hinzugefügt. Weitere Informationen finden Sie unter Deaktivieren Ihrer Pipeline .	7. April 2015

Änderung	Beschreibung	Veröffentlichungsdatum
Vorlagen und Konsole wurden aktualisiert	Neue Vorlagen hinzugefügt. Das Kapitel Erste Schritte wurde aktualisiert und verwendet nun die ShellComm andActivity Vorlage Erste Schritte mit. Weitere Informationen finden Sie unter Erstellen Sie mit der CLI eine Pipeline aus Data Pipeline-Vorlagen.	25. November 2014
VPC-Unterstützung	Unterstützung zum Starten von Ressourcen in einer virtuellen privaten Cloud (VPC) hinzugefügt.	12. März 2014
Regionsunterstützung	Unterstützung für mehrere Serviceregionen hinzugefügt. Darüber us-east-1 hinaus AWS Data Pipeline wird neu-west-1 , ap-northeast-1 ap-southeast-2 , und unterstütztus-west-2 .	20. Februar 2014
Amazon Redshift-Unterstützung	Unterstützung für Amazon Redshift wurde hinzugefügt AWS Data Pipeline, einschließlich einer neuen Konsolenvorlage (Copy to Redshift) und eines Tutorials zur Demonstration der Vorlage. Weitere Informationen dazu finden Sie unter Daten mithilfe von Amazon Redshift nach Amazon Redshift kopieren AWS Data Pipeline , RedshiftDataNode , RedshiftDatabase und RedshiftCopyActivity .	6. November 2013
PigActivity	Hinzugefügt PigActivity, was native Unterstützung für Pig bietet. Weitere Informationen finden Sie unter PigActivity .	15. Oktober 2013
Neue Konsolenvorlage, neue Aktivität und neues Datenformat	Die neue CrossRegion DynamoDB Copy-Konsolenvorlage wurde hinzugefügt, einschließlich der neuen HiveCopyActivity und der Dynamo-Vorlage. DBExport DataFormat	21. August 2013

Änderung	Beschreibung	Veröffentlichungsdatum
Cascading-Ausfälle und erneute Ausführungen	Es wurden Informationen zu AWS Data Pipeline kaskadierenden Fehlern und Wiederholungsverhalten hinzugefügt. Weitere Informationen finden Sie unter Cascading-Ausfälle und erneute Ausführungen .	8. August 2013
Video zur Fehlerbehebung	Das Video zur AWS Data Pipeline grundlegenden Problembehandlung wurde hinzugefügt. Weitere Informationen finden Sie unter Fehlerbehebung .	17. Juli 2013
Bearbeiten von aktiven Pipelines	Zusätzliche Informationen zum Bearbeiten von aktiven Pipelines und erneuten Ausführen von Pipeline-Komponenten hinzugefügt. Weitere Informationen finden Sie unter Bearbeiten Ihrer Pipeline .	17. Juli 2013
Verwenden von Ressourcen in verschiedenen Regionen	Zusätzliche Informationen zum Verwenden von Ressourcen in verschiedenen Regionen hinzugefügt. Weitere Informationen finden Sie unter Verwenden einer Pipeline mit Ressourcen in mehreren Regionen .	17. Juni 2013
Status WAITING_ON_DEPENDENCIES	Status CHECKING_PRECONDITIONS in WAITING_ON_DEPENDENCIES geändert und Laufzeitfeld @waitingOn für Pipeline-Objekte hinzugefügt.	20. Mai 2013
Dynamo-Format DBData	Vorlage im DBData Dynamo-Format hinzugefügt.	23. April 2013
Video zur Verarbeitung von Webprotokollen, Unterstützung von Spot-Instances	Das Video „Webprotokolle mit AWS Data Pipeline, Amazon EMR und Hive verarbeiten“ und die Unterstützung von Amazon EC2 Spot-Instances wurden vorgestellt.	21. Februar 2013
	Die erste Version des AWS Data Pipeline Entwicklerhandbuchs.	20. Dezember 2012