



Guia do usuário do Amazon EMR Sem Servidor

Amazon EMR



Amazon EMR: Guia do usuário do Amazon EMR Sem Servidor

Copyright © 2026 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

As marcas comerciais e imagens de marcas da Amazon não podem ser usadas no contexto de nenhum produto ou serviço que não seja da Amazon, nem de qualquer maneira que possa gerar confusão entre os clientes ou que deprecie ou desprestige a Amazon. Todas as outras marcas comerciais que não pertencem à Amazon pertencem a seus respectivos proprietários, que podem ou não ser afiliados, patrocinados pela Amazon ou ter conexão com ela.

Table of Contents

O que é o Amazon EMR Sem Servidor?	1
Conceitos	1
Versão de lançamento	1
Aplicação	2
Execução de trabalho	3
Operadores	3
Capacidade pré-inicializada	3
EMR Studio	4
Pré-requisitos para começar a usar	5
Inscreva-se para um Conta da AWS	5
Criar um usuário com acesso administrativo	6
Conceder permissões	7
Conceder acesso programático	9
Configure o AWS CLI	11
Abra o console do	12
Introdução	13
Permissões	13
Armazenamento	13
Workloads interativas	13
Criação de um perfil de runtime de trabalhos	14
Como iniciar com o console	19
Etapa 1: Criar um aplicativo do	20
Etapa 2: enviar uma execução de trabalho ou workload interativa	21
Etapa 3: exibir a interface do usuário da aplicação e os logs	24
Etapa 4: limpar	24
Começando a partir do AWS CLI	25
Etapa 1: Criar um aplicativo do	25
Etapa 2: enviar execução de um trabalho	26
Etapa 3: revisar a saída	29
Etapa 4: limpar	30
Como interagir e configurar uma aplicação do EMR Sem Servidor	32
Estados da aplicação	32
Uso do console do EMR Studio	34
Criar uma aplicação do	34

Listagem de aplicações no console do EMR Studio	35
Gerenciamento de aplicações no console do EMR Studio	35
Usando o AWS CLI	36
Como configurar uma aplicação	37
Comportamento da aplicação	37
Capacidade pré-inicializada para trabalhar com uma aplicação no EMR Sem Servidor	39
Configuração padrão de aplicações	43
Personalização de uma imagem	49
Pré-requisitos	38
Etapa 1: criar uma imagem personalizada usando imagens base do EMR Sem Servidor	51
Etapa 2: validar a imagem localmente	51
Etapa 3: upload da imagem em um repositório do Amazon ECR	52
Etapa 4: criar ou atualizar uma aplicação com imagens personalizadas	53
Etapa 5: permitir que o EMR Sem Servidor acesse o repositório de imagens personalizadas	54
Considerações e limitações	55
Configuração do acesso à VPC para que aplicações do EMR Sem Servidor se conectem aos dados	56
Criar aplicação	56
Configuração de aplicações	59
Práticas recomendadas do planejamento de sub-redes	60
Opções de arquitetura	61
Uso da arquitetura x86_64	62
Uso da arquitetura arm64 (Graviton)	62
Lançamento de aplicações com o Graviton	62
Conversão de aplicações existentes para o Graviton	63
Considerações	64
Simultaneidade e enfileiramento de trabalhos	64
Principais benefícios da simultaneidade e do enfileiramento	65
Conceitos básicos de simultaneidade e enfileiramento	65
Considerações sobre simultaneidade e enfileiramento	66
Upload dos dados	68
Pré-requisitos	68
Conceitos básicos da classe S3 Express One Zone	69
Execução de trabalhos	71
Estados de execução de trabalho	71

Cancelamento da execução de trabalhos com período de tolerância	73
Período de tolerância para trabalhos em lote	73
Período de tolerância para trabalhos de transmissão	75
Considerações	55
Uso do console do EMR Studio	78
Enviar uma tarefa	78
Acessar execuções de trabalhos	81
Usando o AWS CLI	81
Política do IAM de execução	82
Conceitos básicos	83
Exemplos de comandos da CLI	83
Observações importantes	85
Interseção de políticas	85
Uso de discos otimizados para embaralhamento	88
Benefícios principais	88
Introdução	89
Usando armazenamento sem servidor	93
Benefícios principais	93
Introdução	94
Considerações e limitações	95
Suportado Regiões da AWS	96
Trabalhos de streaming para processamento contínuo de dados transmitidos	96
Considerações e limitações	98
Introdução	99
Conectores de streaming	100
Gerenciamento de logs	102
Uso das configurações do Spark ao executar trabalhos do EMR Sem Servidor	103
Parâmetros do Spark	103
Propriedades do Spark	107
Práticas recomendadas de configuração de recursos	112
Exemplos do Spark	113
Uso das configurações do Hive ao executar trabalhos do EMR Sem Servidor	114
Parâmetros do Hive	115
Propriedades do Hive	117
Exemplos do Hive	131
Resiliência no trabalho	132

Monitoramento de um trabalho com uma política de repetição	136
Logs com política de nova tentativa	136
Configuração da metastore para EMR Sem Servidor	136
Usando o AWS Glue Data Catalog como metastore	137
Uso de uma metastore externa do Hive	142
Trabalhando com a hierarquia de vários catálogos do AWS Glue no EMR Serverless	147
Considerações sobre o uso de uma metastore externa	148
Acesso entre contas do S3	148
Pré-requisitos	149
Uso de uma política de bucket do S3	149
Uso de um perfil assumido	150
Exemplos de perfis assumidos	153
Solucionar de problemas de erros	157
Erro: o trabalho falhou porque a conta atingiu o limite do serviço para o número máximo de vCPUs que pode usar simultaneamente.	158
Erro: o trabalho falhou porque a aplicação excedeu as configurações de capacidade máxima.	158
Erro: o trabalho falhou porque o operador não pôde ser alocado porque a aplicação excedeu a capacidade máxima.	158
Erro: acesso ao S3 negado. Verifique as permissões de acesso ao S3 do perfil de runtime do trabalho nos recursos necessários do S3.	159
Erro ModuleNotFoundError: Nenhum módulo nomeado<module>. Consulte o guia do usuário sobre como usar bibliotecas Python com o EMR Sem Servidor.	159
Erro: Não foi possível assumir o perfil de execução <role name> porque ele não existe ou não está configurado com o relacionamento de confiança necessário.	159
Alocação de custos em nível de trabalho	159
Comportamento padrão do	160
Como ativar ou desativar o recurso	160
Condições e limitações	161
Execução de workloads interativas	162
Visão geral do	162
Pré-requisitos	162
Permissões	163
Configuração	164
Considerações	164
Execução de workloads interativas por meio do endpoint do Apache Livy	166

Pré-requisitos	166
Permissões obrigatórias	166
Introdução	168
Considerações	175
Registro em log e monitoramento	177
Armazenamento de logs	177
Armazenamento gerenciado	178
Amazon S3	179
Amazon CloudWatch	181
Logs alternados	184
Criptografia de logs	186
Armazenamento gerenciado	186
Buckets do Amazon S3	186
Amazon CloudWatch	186
Permissões obrigatórias	187
Configuração do Log4j2	191
Log4j2 e Spark	191
Monitoramento	195
Aplicações e trabalhos	196
Métricas do mecanismo do Spark	206
Métricas de uso	210
Automatizando com EventBridge	211
Exemplos de eventos do EMR Serverless EventBridge	212
Marcando atributos	217
O que é uma tag?	217
Marcar recursos	217
Limitações de tags	218
Como trabalhar com tags	219
Tutoriais	221
Uso do Java 17	221
JAVA_HOME	221
spark-defaults	222
Uso do Hudi	223
Uso do Iceberg	224
Uso de bibliotecas Python	225
Uso dos recursos nativos do Python	225

Criação de um ambiente virtual Python	226
Configurando PySpark trabalhos para usar bibliotecas Python	227
Uso de diferentes versões do Python	228
Uso do OSS no Delta Lake	230
Amazon EMR 6.9.0 e versões posteriores	230
Amazon EMR 6.8.0 e versões anteriores	231
Envio de trabalhos do Airflow	232
Uso de funções definidas pelo usuário no Hive	235
Uso de imagens personalizada	236
Uso de uma versão personalizada do Python	237
Uso de uma versão personalizada do Java	237
Criação de uma imagem de ciência de dados	238
Processamento de dados geoespaciais com o Apache Sedona	238
Informações de licenciamento para o uso de imagens personalizadas	239
Usar o Spark no Amazon Redshift	240
Iniciar uma aplicação do Spark	240
Autenticação no Amazon Redshift	241
Leitura e gravação para o Amazon Redshift	244
Considerações	246
Conectar-se ao DynamoDB	247
Etapa 1: fazer upload no Amazon S3	247
Etapa 2: criar uma tabela do Hive	248
Etapa 3: copiar para o DynamoDB	249
Etapa 4: consulta do DynamoDB	251
Configurar o acesso entre contas	253
Considerações	255
Segurança	257
Práticas recomendadas de segurança	258
Aplicação do princípio de privilégio mínimo	258
Isolamento de código de uma aplicação não confiável	258
Permissões de controle de acesso por perfil (RBAC)	258
Proteção de dados	259
Criptografia em repouso	260
Criptografia em trânsito	262
Identity and Access Management (IAM)	263
Público	263

Autenticação com identidades	264
Gerenciar o acesso usando políticas	265
Como o EMR Sem Servidor funciona com o IAM	267
Uso de perfis vinculados ao serviço	272
Perfis de runtime do trabalho para o Amazon EMR Sem Servidor	278
Políticas de acesso do usuário	281
Políticas para controle de acesso baseado em etiquetas	285
Políticas baseadas em identidade	289
Atualizações da política	292
Solução de problemas	293
Propagação de identidades confiáveis	295
Visão geral do	295
Recursos e benefícios	296
Como funciona	296
Introdução à propagação de identidade confiável	297
Propagação de identidade confiável para cargas de trabalho interativas	301
Sessões de usuário em segundo plano	301
Considerações sobre a integração sem servidor do Trusted-Identity-Propagation EMR	306
Usar o Lake Formation com o EMR Sem Servidor	307
Disponibilidade de recursos	308
Acesso completo à tabela do Lake Formation para o EMR Serverless	308
Lake Formation para o FGAC	326
Criptografia entre trabalhadores	358
Como habilitar a criptografia TLS mútua no EMR Sem Servidor	359
Criptografia de disco com KMS CMK	359
Usando o contexto de criptografia	360
Configurando a criptografia de disco com chaves gerenciadas pelo cliente	360
Permissões necessárias para criptografia de disco	362
Monitorando o uso da chave	365
Saiba mais	367
Secrets Manager para proteção de dados	368
Como funcionam os segredos	368
Criar um segredo	369
Especificação de referências secretas	369
Concessão de acesso ao segredo	371
Alternância do segredo	373

Concessão de Acesso do S3 para o controle de acesso aos dados	374
Visão geral do	374
Como executar uma aplicação	374
Considerações	376
CloudTrail para registro	376
Informações do EMR Serverless em CloudTrail	376
Noções básicas das entradas do arquivo de log do EMR Sem Servidor	377
Validação de conformidade	379
Resiliência	380
Segurança da infraestrutura	380
Análise de configuração e vulnerabilidade	381
Endpoints e cotas	382
Service endpoints	382
Cotas de serviço	387
Limites de API	388
Outras considerações	55
Versões de liberação	392
AWS runtime for Apache Spark(pré-visualização do emr-spark-8.0)	393
EMR Serverless7.12.0	395
EMR Serverless7.11.0	396
EMR Serverless7.10.0	396
EMR Serverless7.9.0	397
EMR Serverless7.8.0	397
EMR Serverless7.7.0	397
EMR Serverless7.6.0	398
EMR Serverless7.5.0	398
EMR Serverless7.4.0	398
EMR Serverless7.3.0	399
EMR Serverless7.2.0	399
EMR Serverless7.1.0	400
EMR Serverless7.0.0	400
EMR Serverless6.15.0	400
EMR Serverless6.14.0	401
EMR Serverless6.13.0	401
EMR Serverless6.12.0	402
EMR Serverless6.11.0	402

EMR Serverless6.10.0	403
EMR Serverless6.9.0	403
EMR Serverless6.8.0	404
EMR Serverless6.7.0	404
Alterações específicas do mecanismo	405
EMR Serverless6.6.0	405
Histórico do documento	407
.....	cdix

O que é o Amazon EMR Sem Servidor?

O Amazon EMR Sem Servidor é uma opção de implantação do Amazon EMR que fornece um ambiente de runtime com tecnologia sem servidor. Isso simplifica a operação de aplicações de analytics que usam as estruturas de código aberto mais recentes, como o Apache Spark e o Apache Hive. Com o EMR Sem Servidor, você não precisa configurar, otimizar, proteger ou operar clusters para executar aplicações com essas estruturas.

O EMR Sem Servidor ajuda você a evitar o provisionamento excessivo ou insuficiente de recursos em trabalhos de processamento de dados. O EMR Sem Servidor determina automaticamente os recursos de que a aplicação precisa, faz com que esses recursos processem seus trabalhos e os libera quando os trabalhos são concluídos. Em casos de uso em que as aplicações precisam de uma resposta em segundos, como análise interativa de dados, você pode pré-inicializar os recursos de que a aplicação precisa ao criá-la.

Com o EMR Sem Servidor, você continuará a obter os benefícios do Amazon EMR, como compatibilidade de código aberto, simultaneidade e performance de runtime otimizada para estruturas populares.

O EMR Sem Servidor é adequado para clientes que desejam facilidade na operação de aplicações usando estruturas de código aberto. Ele oferece inicialização rápida de trabalhos, gerenciamento automático de capacidade e controles de custos diretos.

Conceitos

Nesta seção, abordamos os termos e conceitos do EMR Sem Servidor que aparecem no Guia do usuário do EMR Sem Servidor.

Versão de lançamento

Uma versão do Amazon EMR corresponde a um conjunto de aplicações de código aberto do ecossistema de big data. Cada versão inclui diferentes aplicações, componentes e recursos de big data que você seleciona para que o EMR Sem Servidor implante e configure, de modo que eles possam executar suas aplicações. Ao criar uma aplicação, especifique sua versão de lançamento. Escolha a versão de lançamento do Amazon EMR e a versão da estrutura de código aberto que deseja usar na aplicação. Para saber mais sobre versões de pré-lançamento, consulte [Versões de lançamento do Amazon EMR Sem Servidor](#).

Aplicação

Com o EMR Sem Servidor, você pode criar uma ou mais aplicações do EMR Sem Servidor que usam estruturas de analytics de código aberto. Para criar uma aplicação, especifique os seguintes atributos:

- A versão de lançamento do Amazon EMR para a versão da estrutura de código aberto que você deseja usar. Para determinar sua versão de lançamento, consulte [Versões de lançamento do Amazon EMR Sem Servidor](#).
- O runtime específico que você deseja que a aplicação use, como o Apache Spark ou o Apache Hive.

Depois de criar uma aplicação, você pode enviar trabalhos de processamento de dados ou solicitações interativas para ela.

Cada aplicação do EMR Sem Servidor é executada em uma Amazon Virtual Private Cloud (VPC) segura, estritamente separada de outras aplicações. Além disso, use políticas AWS Identity and Access Management (IAM) para definir quais usuários e funções podem acessar o aplicativo. Você também pode especificar limites para controlar e rastrear os custos de uso incorridos pela aplicação.

Considere criar várias aplicações quando precisar fazer o seguinte:

- Usa diferentes estruturas de código aberto
- Usar versões diferentes de estruturas de código aberto para diferentes casos de uso
- Execute A/B testes ao atualizar de uma versão para outra
- Manter ambientes lógicos separados para cenários de teste e produção
- Fornecer ambientes lógicos separados para equipes diferentes com controles de custos e rastreamento de uso independentes
- Separe line-of-business aplicativos diferentes

O EMR Sem Servidor é um serviço regional que simplifica a forma como as workloads são executadas em várias zonas de disponibilidade em uma região. Para saber mais sobre como usar aplicações com o EMR Sem Servidor, consulte [Como interagir e configurar uma aplicação do EMR Sem Servidor](#).

Execução de trabalho

A execução de um trabalho é uma solicitação enviada a uma aplicação do EMR Sem Servidor que a aplicação executa e acompanha de forma assíncrona até a conclusão. Exemplos de trabalhos incluem uma consulta HiveQL que você envia para um aplicativo Apache Hive ou um script de processamento de dados que você envia para PySpark um aplicativo Apache Spark. Ao enviar um trabalho, você deve especificar uma função de tempo de execução, criada no IAM, que o trabalho usa para acessar AWS recursos, como objetos do Amazon S3. Você pode enviar várias solicitações de execução de trabalho para um aplicativo, e cada execução de trabalho pode usar uma função de tempo de execução diferente para acessar AWS recursos. Uma aplicação do EMR Sem Servidor começa a executar trabalhos assim que os recebe e executa várias solicitações de trabalho simultaneamente. Para saber mais sobre como o EMR Sem Servidor executa trabalhos, consulte [Execução de trabalhos](#).

Operadores

Uma aplicação do EMR Sem Servidor usa trabalhadores internamente para executar workloads. Os tamanhos padrão desses trabalhadores são baseados no tipo de aplicação e na versão de lançamento do Amazon EMR. Ao programar uma execução de trabalho, substitua esses tamanhos.

Quando você envia um trabalho, o EMR Sem Servidor calcula os recursos que a aplicação precisa para o trabalho e agenda os trabalhadores. O EMR Sem Servidor divide as workloads em tarefas, baixa imagens, provisiona e configura trabalhadores e os desativa quando o trabalho é concluído. O EMR Sem Servidor escala automaticamente os trabalhadores com base na workload e no paralelismo necessários em cada estágio do trabalho. Esse ajuste de escala automático elimina a necessidade de estimar o número de trabalhadores que a aplicação precisa para executar as workloads.

Capacidade pré-inicializada

O EMR Sem Servidor fornece um recurso de capacidade pré-inicializada que mantém os trabalhadores inicializados e prontos para responder em segundos. Essa capacidade cria efetivamente um grupo de aquecimento de trabalhadores para uma aplicação. Para configurar esse recurso para cada aplicação, defina o parâmetro `initial-capacity` de uma aplicação. Quando você configura a capacidade pré-inicializada, os trabalhos podem começar imediatamente para que você possa implementar aplicações iterativas e trabalhos urgentes. Para saber mais sobre trabalhadores pré-inicializados, consulte [Configuração de uma aplicação ao trabalhar com o EMR Sem Servidor](#).

EMR Studio

O EMR Studio é o console do usuário para gerenciar suas aplicações do EMR Sem Servidor. Se não existir um EMR Studio em sua conta quando criar sua primeira aplicação do EMR Sem Servidor, criaremos um para você automaticamente. Acesse o EMR Studio no console do Amazon EMR ou ative o acesso federado do seu provedor de identidades (IdP) por meio do IAM ou do Centro de Identidade do IAM. Ao fazer isso, os usuários podem acessar o Studio e gerenciar aplicações do EMR Sem Servidor sem acesso direto ao console do Amazon EMR. Para saber mais sobre como as aplicações do EMR Sem Servidor funcionam com o EMR Studio, consulte [Criação de uma aplicação do EMR Sem Servidor usando o console do EMR Studio](#) e [Execução de trabalhos no console do EMR Studio](#).

Pré-requisitos para começar a usar o EMR Sem Servidor

Esta seção descreve os pré-requisitos administrativos para executar o EMR Sem Servidor. Isso inclui configuração da conta e gerenciamento de permissões.

Tópicos

- [Inscreva-se para um Conta da AWS](#)
- [Criar um usuário com acesso administrativo](#)
- [Conceder permissões](#)
- [Instale e configure o AWS CLI](#)
- [Abra o console do .](#)

Inscreva-se para um Conta da AWS

Se você não tiver um Conta da AWS, conclua as etapas a seguir para criar um.

Para se inscrever em um Conta da AWS

1. Abra a <https://portal.aws.amazon.com/billing/inscrição>.
2. Siga as instruções online.

Parte do procedimento de inscrição envolve receber uma chamada telefônica ou uma mensagem de texto e inserir um código de verificação pelo teclado do telefone.

Quando você se inscreve em um Conta da AWS, um Usuário raiz da conta da AWS é criado. O usuário-raiz tem acesso a todos os Serviços da AWS e recursos na conta. Como prática recomendada de segurança, atribua o acesso administrativo a um usuário e use somente o usuário-raiz para executar [tarefas que exigem acesso de usuário-raiz](#).

AWS envia um e-mail de confirmação após a conclusão do processo de inscrição. A qualquer momento, você pode visualizar a atividade atual da sua conta e gerenciar sua conta acessando <https://aws.amazon.com/e> escolhendo Minha conta.

Criar um usuário com acesso administrativo

Depois de se inscrever em um Conta da AWS, proteja seu Usuário raiz da conta da AWS AWS IAM Identity Center, habilite e crie um usuário administrativo para que você não use o usuário root nas tarefas diárias.

Proteja seu Usuário raiz da conta da AWS

1. Faça login [Console de gerenciamento da AWS](#) como proprietário da conta escolhendo Usuário raiz e inserindo seu endereço de Conta da AWS e-mail. Na próxima página, insira a senha.

Para obter ajuda ao fazer login usando o usuário-raiz, consulte [Fazer login como usuário-raiz](#) no Guia do usuário do Início de Sessão da AWS .

2. Habilite a autenticação multifator (MFA) para o usuário-raiz.

Para obter instruções, consulte [Habilitar um dispositivo de MFA virtual para seu usuário Conta da AWS raiz \(console\) no Guia](#) do usuário do IAM.

Criar um usuário com acesso administrativo

1. Habilita o Centro de Identidade do IAM.

Para obter instruções, consulte [Habilitar o AWS IAM Identity Center](#) no Guia do usuário do AWS IAM Identity Center .

2. No Centro de Identidade do IAM, conceda o acesso administrativo a um usuário.

Para ver um tutorial sobre como usar o Diretório do Centro de Identidade do IAM como fonte de identidade, consulte [Configurar o acesso do usuário com o padrão Diretório do Centro de Identidade do IAM](#) no Guia AWS IAM Identity Center do usuário.

Iniciar sessão como o usuário com acesso administrativo

- Para fazer login com o seu usuário do Centro de Identidade do IAM, use o URL de login enviado ao seu endereço de e-mail quando o usuário do Centro de Identidade do IAM foi criado.

Para obter ajuda para fazer login usando um usuário do IAM Identity Center, consulte [Como fazer login no portal de AWS acesso](#) no Guia Início de Sessão da AWS do usuário.

Atribuir acesso a usuários adicionais

1. No Centro de Identidade do IAM, crie um conjunto de permissões que siga as práticas recomendadas de aplicação de permissões com privilégio mínimo.

Para obter instruções, consulte [Criar um conjunto de permissões](#) no Guia do usuário do AWS IAM Identity Center .

2. Atribua usuários a um grupo e, em seguida, atribua o acesso de logon único ao grupo.

Para obter instruções, consulte [Adicionar grupos](#) no Guia do usuário do AWS IAM Identity Center .

Conceder permissões

Em ambientes de produção, sugerimos que você use políticas mais refinadas. Para ver alguns exemplos dessas políticas, consulte [Exemplos de políticas de acesso de usuários do EMR Sem Servidor](#). Para saber mais sobre gerenciamento de acesso, consulte [Gerenciamento de acesso para AWS recursos](#) no Guia do usuário do IAM.

Para usuários que precisam começar a usar o EMR Sem Servidor em um ambiente sandbox, use uma política semelhante à seguinte:

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "EMRStudioCreate",
      "Effect": "Allow",
      "Action": [
        "elasticmapreduce:CreateStudioPresignedUrl",
        "elasticmapreduce:DescribeStudio",
        "elasticmapreduce:CreateStudio",
        "elasticmapreduce:ListStudios"
      ],
      "Resource": [
        "*"
      ]
    }
  ],
}
```

```
{
  "Sid": "EMRServerlessFullAccess",
  "Effect": "Allow",
  "Action": [
    "emr-serverless:*"
  ],
  "Resource": [
    "*"
  ]
},
{
  "Sid": "AllowEC2ENICreationWithEMRTags",
  "Effect": "Allow",
  "Action": [
    "ec2:CreateNetworkInterface"
  ],
  "Resource": [
    "arn:aws:ec2:*:*:network-interface/*"
  ],
  "Condition": {
    "StringEquals": {
      "aws:CalledViaLast": "ops.emr-serverless.amazonaws.com"
    }
  }
},
{
  "Sid": "AllowEMRServerlessServiceLinkedRoleCreation",
  "Effect": "Allow",
  "Action": [
    "iam:CreateServiceLinkedRole"
  ],
  "Resource": [
    "arn:aws:iam::*:role/aws-service-role/*"
  ]
}
]
```

Para conceder acesso, adicione as permissões aos seus usuários, grupos ou perfis:

- Usuários e grupos em AWS IAM Identity Center:

Crie um conjunto de permissões. Siga as instruções em [Criação de um conjunto de permissões](#) no Guia do usuário do AWS IAM Identity Center .

- Usuários gerenciados no IAM com provedor de identidades:

Crie um perfil para a federação de identidades. Siga as instruções em [Criando um perfil para um provedor de identidades de terceiros \(federação\)](#) no Guia do Usuário do IAM.

- Usuários do IAM:
 - Crie um perfil que seu usuário possa assumir. Siga as instruções em [Criação de um perfil para um usuário do IAM](#) no Guia do usuário do IAM.
 - (Não recomendado) Vincule uma política diretamente a um usuário ou adicione um usuário a um grupo de usuários. Siga as instruções em [Adição de permissões a um usuário \(console\)](#) no Guia do usuário do IAM.

Conceder acesso programático

Os usuários precisam de acesso programático se quiserem interagir com pessoas AWS fora do Console de gerenciamento da AWS. A forma de conceder acesso programático depende do tipo de usuário que está acessando AWS.

Para conceder acesso programático aos usuários, selecione uma das seguintes opções:

Qual usuário precisa de acesso programático?	Para	Por
IAM	(Recomendado) Use as credenciais do console como credenciais temporárias para assinar solicitações programáticas para o AWS CLI, AWS SDKs ou. AWS APIs	<p>Siga as instruções da interface que deseja utilizar.</p> <ul style="list-style-type: none"> • Para o AWS CLI, consulte Login para desenvolvimento AWS local no Guia AWS Command Line Interface do usuário. • Para AWS SDKs isso, consulte Login para desenvolvimento AWS local no Guia de referência de

Qual usuário precisa de acesso programático?	Para	Por
<p>Identidade da força de trabalho</p> <p>(Usuários gerenciados no Centro de Identidade do IAM)</p>	<p>Use credenciais temporárias para assinar solicitações programáticas para o AWS CLI AWS SDKs, ou. AWS APIs</p>	<p>ferramentas AWS SDKs e ferramentas.</p> <p>Siga as instruções da interface que deseja utilizar.</p> <ul style="list-style-type: none"> • Para o AWS CLI, consulte Configurando o AWS CLI para uso AWS IAM Identity Center no Guia do AWS Command Line Interface usuário. • Para AWS SDKs, ferramentas e AWS APIs, consulte a autenticação do IAM Identity Center no Guia de referência de ferramentas AWS SDKs e ferramentas.
<p>IAM</p>	<p>Use credenciais temporárias para assinar solicitações programáticas para o AWS CLI AWS SDKs, ou. AWS APIs</p>	<p>Siga as instruções em Como usar credenciais temporárias com AWS recursos no Guia do usuário do IAM.</p>

Qual usuário precisa de acesso programático?	Para	Por
IAM	(Não recomendado) Use credenciais de longo prazo para assinar solicitações programáticas para o AWS CLI, AWS SDKs, ou. AWS APIs	<p>Siga as instruções da interface que deseja utilizar.</p> <ul style="list-style-type: none"> • Para isso AWS CLI, consulte Autenticação usando credenciais de usuário do IAM no Guia do AWS Command Line Interface usuário. • Para ferramentas AWS SDKs e ferramentas, consulte Autenticar usando credenciais de longo prazo no Guia de referência de ferramentas AWS SDKs e ferramentas. • Para isso AWS APIs, consulte Gerenciamento de chaves de acesso para usuários do IAM no Guia do usuário do IAM.

Instale e configure o AWS CLI

Se você quiser usar o EMR Serverless APIs, instale a versão mais recente do (). AWS Command Line Interface AWS CLI Você não AWS CLI precisa usar o EMR Serverless do console do EMR Studio e começar sem a CLI seguindo as etapas em. [Conceitos básicos do EMR Sem Servidor usando o console](#)

Para configurar o AWS CLI

1. Para instalar a versão mais recente do AWS CLI para macOS, Linux ou Windows, consulte [Instalando ou atualizando a versão mais recente do. AWS CLI](#)

2. [Para configurar AWS CLI e proteger seu acesso Serviços da AWS, incluindo o EMR Serverless, consulte Configuração rápida com. aws configure](#)
3. Para verificar a configuração, digite o seguinte DataBrew comando no prompt de comando.

```
aws emr-serverless help
```

AWS CLI os comandos usam o padrão Região da AWS da sua configuração, a menos que você o defina com um parâmetro ou um perfil. Para definir seu Região da AWS com um parâmetro, adicione o `--region` parâmetro a cada comando.

Para definir seu Região da AWS com um perfil, primeiro adicione um perfil nomeado no `~/.aws/config` arquivo ou no `%UserProfile%/.aws/config` arquivo (para Microsoft Windows). Siga as etapas em [Perfis nomeados para a AWS CLI](#). Em seguida, defina suas Região da AWS e outras configurações com um comando semelhante ao do exemplo a seguir.

```
[profile emr-serverless]
aws_access_key_id = ACCESS-KEY-ID-OF-IAM-USER
aws_secret_access_key = SECRET-ACCESS-KEY-ID-OF-IAM-USER
region = us-east-1
output = text
```

Abra o console do .

A maioria dos tópicos orientados para o console nesta seção começa no [console do Amazon EMR](#). Se você ainda não estiver conectado ao seu Conta da AWS, faça login, abra o [console do Amazon EMR](#) e vá para a próxima seção para continuar começando a usar o Amazon EMR.

Conceitos básicos do Amazon EMR Sem Servidor

Este tutorial ajuda você nos conceitos básicos do EMR Sem Servidor ao implantar um exemplo de workload do Spark ou do Hive. Você criará, executará e depurará a própria aplicação. Mostramos as opções padrão na maior parte deste tutorial.

Antes de iniciar uma aplicação do EMR Sem Servidor, conclua as tarefas a seguir.

Tópicos

- [Concessão de permissões para usar o EMR Sem Servidor](#)
- [Preparação do armazenamento do EMR Sem Servidor](#)
- [Criação de um EMR Studio para executar workloads interativas](#)
- [Criação de um perfil de runtime de trabalhos](#)
- [Conceitos básicos do EMR Sem Servidor usando o console](#)
- [Começando a partir do AWS CLI](#)

Concessão de permissões para usar o EMR Sem Servidor

Para usar o EMR Sem Servidor, você precisa de um usuário ou perfil do IAM com uma política anexada que conceda permissões para o EMR Sem Servidor. Para criar um usuário e anexar a política apropriada a ele, siga as instruções em [Conceder permissões](#).

Preparação do armazenamento do EMR Sem Servidor

Neste tutorial, você usará um bucket do S3 para armazenar arquivos e logs de saída do exemplo de workload do Spark ou do Hive que será executada usando uma aplicação do EMR Sem Servidor. Para criar um bucket, siga as instruções em [Criação de um bucket](#) no Guia do usuário do console do Amazon Simple Storage Service. Substitua qualquer outra referência a *amzn-s3-demo-bucket* pelo nome do bucket recém-criado.

Criação de um EMR Studio para executar workloads interativas

Se você quiser usar o EMR Sem Servidor para executar consultas interativas por meio de cadernos hospedados no EMR Studio, será necessário especificar um bucket do S3 e o [perfil de serviço mínimo para o EMR Sem Servidor](#) criar um Workspace. Para conhecer as etapas de configuração,

consulte [Set up an EMR Studio](#) no Guia de gerenciamento do Amazon EMR. Para obter mais informações sobre workloads interativas, consulte [Execução de workloads interativas com o EMR Sem Servidor por meio do EMR Studio](#).

Criação de um perfil de runtime de trabalhos

Os trabalhos executados no EMR Serverless usam uma função de tempo de execução que fornece permissões granulares para recursos específicos Serviços da AWS e em tempo de execução. Neste tutorial, um bucket público do S3 hospeda os dados e os scripts. O bucket *amzn-s3-demo-bucket* armazena a saída.

Para configurar um perfil de runtime de trabalhos, primeiro crie um perfil de runtime com uma política de confiança para que o EMR Sem Servidor possa usar o novo perfil. Em seguida, anexe a política de acesso do S3 necessária a esse perfil. As etapas a seguir guiam você pelo processo.

Console

1. Navegue até o console do IAM em <https://console.aws.amazon.com/iam/>.
2. No painel de navegação à esquerda, escolha Políticas.
3. Escolha Create Policy.
4. A página Criar política é aberta em uma nova guia. Selecione o Editor de políticas como Json e cole o JSON da política abaixo.

Important

Substitua *amzn-s3-demo-bucket* na política abaixo pelo nome real do bucket criado em [Preparação do armazenamento do EMR Sem Servidor](#). Essa é uma política básica para acesso ao S3. Para obter mais exemplos de perfis de runtime de trabalhos, consulte [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#).

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "ReadAccessForEMRSamples",
```

```

    "Effect": "Allow",
    "Action": [
      "s3:GetObject",
      "s3:ListBucket"
    ],
    "Resource": [
      "arn:aws:s3::*.elasticmapreduce",
      "arn:aws:s3::*.elasticmapreduce/*"
    ]
  },
  {
    "Sid": "FullAccessToOutputBucket",
    "Effect": "Allow",
    "Action": [
      "s3:PutObject",
      "s3:GetObject",
      "s3:ListBucket",
      "s3:DeleteObject"
    ],
    "Resource": [
      "arn:aws:s3:::amzn-s3-demo-bucket",
      "arn:aws:s3:::amzn-s3-demo-bucket/*"
    ]
  },
  {
    "Sid": "GlueCreateAndReadDataCatalog",
    "Effect": "Allow",
    "Action": [
      "glue:GetDatabase",
      "glue:CreateDatabase",
      "glue:GetDataBases",
      "glue:CreateTable",
      "glue:GetTable",
      "glue:UpdateTable",
      "glue>DeleteTable",
      "glue:GetTables",
      "glue:GetPartition",
      "glue:GetPartitions",
      "glue:CreatePartition",
      "glue:BatchCreatePartition",
      "glue:GetUserDefinedFunctions"
    ],
    "Resource": [
      "*"
    ]
  }

```

```

    ]
  }
]
}

```

5. Selecione Avançar para inserir um nome para sua política, como `EMRServerlessS3AndGlueAccessPolicy` e Criar política
6. No painel de navegação esquerdo do console do IAM, escolha Perfis.
7. Selecione Criar perfil.
8. Em tipo de perfil, escolha Política de confiança personalizada e cole a política de confiança a seguir. Isso permite que trabalhos enviados para seus aplicativos Amazon EMR Serverless acessem outros Serviços da AWS em seu nome.

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "sts:AssumeRole"
      ],
      "Resource": "arn:aws:iam::123456789012:role/EMRServerlessExecutionRole",
      "Sid": "AllowSTSAssumerole"
    }
  ]
}

```

9. Escolha Avançar para navegar até a página Adicionar permissões e, em seguida, escolha `EMRServerlessS3 AndGlueAccessPolicy`.
10. Na página Nomear, revisar e criar, em Nome do perfil, digite um nome para o perfil, por exemplo, `EMRServerlessS3RuntimeRole`. Para criar esse perfil do IAM, escolha Criar perfil.

CLI

1. Crie um arquivo denominado `emr-serverless-trust-policy.json` que contenha a política de confiança a ser usada para a função do IAM. O arquivo deve ter a política a seguir.

JSON


```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "EMRServerlessTrustPolicy",
      "Action": [
        "sts:AssumeRole"
      ],
      "Effect": "Allow",
      "Resource": "arn:aws:iam::123456789012:role/EMRServerlessExecutionRole"
    }
  ]
}
```

2. Crie um perfil do IAM denominado `EMRServerlessS3RuntimeRole`. Use a política de confiança criada na etapa anterior.

```
aws iam create-role \
  --role-name EMRServerlessS3RuntimeRole \
  --assume-role-policy-document file://emr-serverless-trust-policy.json
```

Anote o ARN na saída. Você usa o ARN do novo perfil durante o envio do trabalho, depois chamado de *job-role-arn*.

3. Crie um arquivo chamado `emr-sample-access-policy.json` que defina a política do IAM para sua workload. Isso fornece acesso de leitura ao script e aos dados armazenados em buckets públicos do S3 e acesso de leitura e gravação a *amzn-s3-demo-bucket*.

 Important

Substitua *amzn-s3-demo-bucket* na política abaixo pelo nome real do bucket criado em [Preparação do armazenamento do EMR Sem Servidor](#). Essa é uma

política básica para acesso ao AWS Glue e ao S3. Para obter mais exemplos de perfis de runtime de trabalhos, consulte [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#).

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "ReadAccessForEMRSamples",
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3::*.elasticmapreduce",
        "arn:aws:s3::*.elasticmapreduce/*"
      ]
    },
    {
      "Sid": "FullAccessToOutputBucket",
      "Effect": "Allow",
      "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "s3:ListBucket",
        "s3:DeleteObject"
      ],
      "Resource": [
        "arn:aws:s3::amzn-s3-demo-bucket",
        "arn:aws:s3::amzn-s3-demo-bucket/*"
      ]
    },
    {
      "Sid": "GlueCreateAndReadDataCatalog",
      "Effect": "Allow",
      "Action": [
        "glue:GetDatabase",
        "glue:CreateDatabase",

```

```

    "glue:GetDataBases",
    "glue:CreateTable",
    "glue:GetTable",
    "glue:UpdateTable",
    "glue>DeleteTable",
    "glue:GetTables",
    "glue:GetPartition",
    "glue:GetPartitions",
    "glue:CreatePartition",
    "glue:BatchCreatePartition",
    "glue:GetUserDefinedFunctions"
  ],
  "Resource": [
    "*"
  ]
}
]
}

```

4. Crie uma política do IAM chamada `EMRServerlessS3AndGlueAccessPolicy` com o arquivo de política criado na Etapa 3. Anote o ARN na saída, pois você usará o ARN da nova política na próxima etapa.

```

aws iam create-policy \
  --policy-name EMRServerlessS3AndGlueAccessPolicy \
  --policy-document file://emr-sample-access-policy.json

```

Anote o ARN da nova política na saída. Você o substituirá por *policy-arn* na próxima etapa.

5. Anexe a política do IAM `EMRServerlessS3AndGlueAccessPolicy` ao perfil de runtime do trabalho `EMRServerlessS3RuntimeRole`.

```

aws iam attach-role-policy \
  --role-name EMRServerlessS3RuntimeRole \
  --policy-arn policy-arn

```

Conceitos básicos do EMR Sem Servidor usando o console

Esta seção descreve o trabalho com o EMR Sem Servidor, incluindo a criação de um EMR Studio. Também descreve como enviar execuções de trabalhos e exibir logs.

Etapas a serem executadas

- [Etapa 1: criar uma aplicação do EMR Sem Servidor](#)
- [Etapa 2: enviar uma execução de trabalho ou workload interativa](#)
- [Etapa 3: exibir a interface do usuário da aplicação e os logs](#)
- [Etapa 4: limpar](#)

Etapa 1: criar uma aplicação do EMR Sem Servidor

Crie uma aplicação com o EMR Sem Servidor da maneira a seguir.

1. [Faça login no Console de gerenciamento da AWS e abra o console do Amazon EMR em https://console.aws.amazon.com/emr.](https://console.aws.amazon.com/emr)
2. No painel de navegação à esquerda, escolha EMR Sem Servidor para navegar até a página de destino do EMR Sem Servidor.
3. Para criar ou gerenciar aplicações do EMR Sem Servidor, você precisa da interface do EMR Studio.
 - Se você já tiver um EMR Studio no Região da AWS local em que deseja criar um aplicativo, selecione Gerenciar aplicativos para navegar até seu EMR Studio ou selecione o estúdio que deseja usar.
 - Se você não tiver um EMR Studio no local em Região da AWS que deseja criar um aplicativo, escolha Começar e, em seguida, escolha Criar e iniciar o Studio. O EMR Sem Servidor cria um EMR Studio para você, para que consiga criar e gerenciar aplicações.
4. Na interface do usuário de Criar Studio que se abre em uma nova guia, insira o nome, o tipo e a versão de lançamento da aplicação. Se você quiser executar somente trabalhos em lotes, selecione Usar configurações padrão somente para trabalhos em lotes. Para workloads interativas, selecione Usar configurações padrão para workloads interativas. Você também pode executar trabalhos em lotes em aplicações interativas com essa opção. Se necessário, você poderá alterar essas configurações posteriormente.

Para obter mais informações, consulte [Create a studio](#).
5. Selecione Criar aplicação para criar sua primeira aplicação.

Continue na próxima seção [Etapa 2: enviar uma execução de trabalho ou workload interativa](#) para enviar uma execução de trabalho ou uma workload interativa.

Etapa 2: enviar uma execução de trabalho ou workload interativa

Spark job run

Neste tutorial, usamos um PySpark script para calcular o número de ocorrências de palavras únicas em vários arquivos de texto. Um bucket do S3 público e somente leitura armazena o script e o conjunto de dados.

Para executar um trabalho do Spark

1. Faça o upload do exemplo de script `wordcount.py` para o novo bucket com o comando a seguir.

```
aws s3 cp s3://us-east-1.elasticmapreduce/emr-containers/samples/wordcount/scripts/wordcount.py s3://amzn-s3-demo-bucket/scripts/
```

2. A conclusão de [Etapa 1: criar uma aplicação do EMR Sem Servidor](#) leva você à página Detalhes da aplicação no EMR Studio. Lá, escolha a opção Enviar trabalho.
3. Na página Enviar trabalho, conclua o procedimento a seguir.
 - No campo Nome, insira o nome que você deseja chamar para a execução do trabalho.
 - No campo Perfil de runtime, digite o nome do perfil criado em [Criação de um perfil de runtime de trabalhos](#).
 - No campo Localização do script, insira `s3://amzn-s3-demo-bucket/scripts/wordcount.py` como URI do S3.
 - No campo Argumentos do script, insira `["s3://amzn-s3-demo-bucket/emr-serverless-spark/output"]`.
 - Na seção Propriedades do Spark, escolha Editar como texto e insira as configurações a seguir.

```
--conf spark.executor.cores=1 --conf spark.executor.memory=4g --  
conf spark.driver.cores=1 --conf spark.driver.memory=4g --conf  
spark.executor.instances=1
```

4. Para iniciar a execução do trabalho, escolha Enviar trabalho.
5. Na guia Execuções de trabalhos, será exibido seu novo trabalho sendo executado com o status Em execução.

Hive job run

Nesta parte do tutorial, criamos uma tabela, inserimos alguns registros e executamos uma consulta de agregação de contagem. Para executar o trabalho do Hive, primeiro crie um arquivo que contenha todas as consultas do Hive a serem executadas como parte de um único trabalho, faça upload do arquivo no S3 e especifique esse caminho do S3 ao iniciar o trabalho do Hive.

Para executar um trabalho do Hive

1. Crie um arquivo chamado `hive-query.q1` que contenha todas as consultas que você deseja executar no trabalho do Hive.

```
create database if not exists emrserverless;
use emrserverless;
create table if not exists test_table(id int);
drop table if exists Values__Tmp__Table__1;
insert into test_table values (1),(2),(2),(3),(3),(3);
select id, count(id) from test_table group by id order by id desc;
```

2. Faça upload de `hive-query.q1` no bucket do S3 com o comando a seguir.

```
aws s3 cp hive-query.q1 s3://amzn-s3-demo-bucket/emr-serverless-hive/query/hive-query.q1
```

3. A conclusão de [Etapa 1: criar uma aplicação do EMR Sem Servidor](#) leva você à página Detalhes da aplicação no EMR Studio. Lá, escolha a opção Enviar trabalho.
4. Na página Enviar trabalho, conclua o procedimento a seguir.
 - No campo Nome, insira o nome que você deseja chamar para a execução do trabalho.
 - No campo Perfil de runtime, digite o nome do perfil criado em [Criação de um perfil de runtime de trabalhos](#).
 - No campo Localização do script, insira `s3://amzn-s3-demo-bucket/emr-serverless-hive/query/hive-query.q1` como URI do S3.
 - Na seção Propriedades do Hive, escolha Editar como texto e insira as configurações a seguir.

```
--hiveconf hive.log.explain.output=false
```

- Na seção Configuração do trabalho, escolha Editar como JSON e insira o JSON a seguir.

```
{
  "applicationConfiguration":
  [{
    "classification": "hive-site",
    "properties": {
      "hive.exec.scratchdir": "s3://amzn-s3-demo-bucket/emr-
serverless-hive/hive/scratch",
      "hive.metastore.warehouse.dir": "s3://amzn-s3-demo-bucket/emr-
serverless-hive/hive/warehouse",
      "hive.driver.cores": "2",
      "hive.driver.memory": "4g",
      "hive.tez.container.size": "4096",
      "hive.tez.cpu.vcores": "1"
    }
  ]
}
```

5. Para iniciar a execução do trabalho, escolha Enviar trabalho.
6. Na guia Execuções de trabalhos, será exibido seu novo trabalho sendo executado com o status Em execução.

Interactive workload

Com o Amazon EMR 6.14.0 e superior, você pode usar cadernos hospedados no EMR Studio para executar workloads interativas para o Spark no EMR Sem Servidor. Para obter mais informações, incluindo permissões e pré-requisitos, consulte [Execução de workloads interativas com o EMR Sem Servidor por meio do EMR Studio](#).

Depois de criar a aplicação e configurar as permissões necessárias, use as seguintes etapas para executar um caderno interativo com o EMR Studio:

1. Navegue até a guia Workspaces no EMR Studio. Se você ainda precisar configurar um local de armazenamento do Amazon S3 e um [perfil de serviço do EMR Studio](#), selecione o botão Configurar Studio no banner na parte superior da tela.
2. Para acessar um caderno, selecione um Workspace ou crie um. Use o Início rápido para abrir seu Workspace em uma nova guia.
3. Vá para a guia recém-aberta. Selecione o ícone Computação na navegação esquerda. Selecione EMR Sem Servidor como o Tipo de computação.
4. Selecione a aplicação interativa que você criou na seção anterior.

5. No campo Perfil de runtime, insira o nome do perfil do IAM que a aplicação do EMR Sem Servidor pode assumir para a execução do trabalho. Para saber mais sobre os perfis de runtime, consulte [Job runtime roles](#) no Guia do usuário do Amazon EMR Sem Servidor.
6. Selecione Anexar. Esse processo pode levar até um minuto. A página será atualizada quando anexada.
7. Escolha um kernel e inicie um caderno. Você também pode procurar exemplos de cadernos no EMR Sem Servidor e copiá-los para o Workspace. Para acessar os exemplos de cadernos, navegue até o menu `{...}` na navegação à esquerda e acesse os cadernos que têm `serverless` no nome do arquivo do caderno.
8. No caderno, você pode acessar o link do log do driver e um link para a interface do usuário do Apache Spark, uma interface em tempo real que fornece métricas para monitorar seu trabalho. Para obter mais informações, consulte [Monitoring EMR Sem Servidor applications and jobs](#) no Guia do usuário do Amazon EMR Sem Servidor.

Quando você anexa uma aplicação a um Workspace do Studio, o início dela é acionado automaticamente se ainda não estiver em execução. Você também pode pré-iniciar a aplicação e mantê-la pronta antes de anexá-la ao Workspace.

Etapa 3: exibir a interface do usuário da aplicação e os logs

Para exibir a interface do usuário da aplicação, primeiro identifique a execução do trabalho. Uma opção para a interface do usuário do Spark ou a interface do usuário do Hive Tez está disponível na primeira linha de opções para a execução desse trabalho, com base no tipo de trabalho. Selecione a opção apropriada.

Se você escolheu a interface do usuário do Spark, selecione a guia Executores para exibir os logs do driver e dos executores. Se você escolheu a interface do usuário do Hive Tez, selecione a guia Todas as tarefas para exibir os logs.

Depois que o status de execução do trabalho for exibido como Êxito, você poderá exibir a saída do trabalho no bucket do S3.

Etapa 4: limpar

Embora a aplicação criada deva parar automaticamente após 15 minutos de inatividade, ainda recomendamos que você libere recursos que não pretende usar novamente.

Para excluir a aplicação, navegue até a página Listar aplicações. Selecione a aplicação que você criou e escolha Ações → Interromper para interromper a aplicação. Depois que a aplicação estiver no estado STOPPED, selecione a mesma aplicação e escolha Ações → Excluir.

Para obter mais exemplos de execução de trabalhos do Spark e do Hive, consulte [Uso das configurações do Spark ao executar trabalhos do EMR Sem Servidor](#) e [Uso das configurações do Hive ao executar trabalhos do EMR Sem Servidor](#).

Começando a partir do AWS CLI

Comece a usar o EMR Serverless a partir dos comandos AWS CLI with para criar um aplicativo, executar trabalhos, verificar a saída da execução do trabalho e excluir seus recursos.

Etapa 1: criar uma aplicação do EMR Sem Servidor

Use o comando [emr-serverless create-application](#) para criar sua primeira aplicação do EMR Sem Servidor. Você precisa especificar o tipo de aplicação e o rótulo de lançamento do Amazon EMR associado à versão da aplicação que deseja usar. O nome da aplicação é opcional.

Spark

Para criar uma aplicação do Spark, execute o comando a seguir.

```
aws emr-serverless create-application \  
  --release-label emr-6.6.0 \  
  --type "SPARK" \  
  --name my-application
```

Hive

Para criar uma aplicação do Hive, execute o comando a seguir.

```
aws emr-serverless create-application \  
  --release-label emr-6.6.0 \  
  --type "HIVE" \  
  --name my-application
```

Anote o ID da aplicação retornado na saída. Você usará o ID para iniciar a aplicação e durante o envio do trabalho, depois chamado de *application-id*.

Antes de prosseguir para [Etapa 2: enviar uma execução de trabalho à aplicação do EMR Sem Servidor](#), certifique-se de que a aplicação tenha atingido o estado CREATED com a API [get-application](#).

```
aws emr-serverless get-application \  
  --application-id application-id
```

O EMR Sem Servidor cria trabalhadores para acomodar seus trabalhos solicitados. Por padrão, eles são criados sob demanda, mas você também pode especificar uma capacidade pré-inicializada definindo o parâmetro `initialCapacity` ao criar a aplicação. Você também pode limitar a capacidade máxima total que uma aplicação pode usar com o parâmetro `maximumCapacity`. Para saber mais sobre essas opções, consulte [Configuração de uma aplicação ao trabalhar com o EMR Sem Servidor](#).

Etapa 2: enviar uma execução de trabalho à aplicação do EMR Sem Servidor

Agora, a aplicação do EMR Sem Servidor está pronta para executar trabalhos.

Spark

Nesta etapa, usamos um PySpark script para calcular o número de ocorrências de palavras exclusivas em vários arquivos de texto. Um bucket do S3 público e somente leitura armazena o script e o conjunto de dados. A aplicação envia o arquivo de saída e os dados de log do runtime do Spark aos diretórios de `/output` e `/logs` no bucket do S3 que você criou.

Para executar um trabalho do Spark

1. Use o comando a seguir para copiar o exemplo de script que executaremos no novo bucket.

```
aws s3 cp s3://us-east-1.elasticmapreduce/emr-containers/samples/wordcount/  
scripts/wordcount.py s3://amzn-s3-demo-bucket/scripts/
```

2. No comando a seguir, substitua `application-id` pelo ID da aplicação. Substitua `job-role-arn` pelo ARN do perfil de runtime que você criou em [Criação de um perfil de runtime de trabalhos](#). Substitua `job-run-name` pelo nome que deseja chamar a execução do trabalho. Substitua todas as strings `amzn-s3-demo-bucket` pelo bucket do Amazon S3 que você criou e adicione `/output` ao caminho. Isso cria uma pasta no seu bucket na qual o EMR Sem Servidor pode copiar os arquivos de saída da aplicação.

```
aws emr-serverless start-job-run \  
  --application-id application-id \  
  --execution-role-arn job-role-arn \  
  --name job-run-name \  
  --job-driver '{  
    "sparkSubmit": {  
      "entryPoint": "s3://amzn-s3-demo-bucket/scripts/wordcount.py",  
      "entryPointArguments": ["s3://amzn-s3-demo-bucket/emr-serverless-  
spark/output"],  
      "sparkSubmitParameters": "--conf spark.executor.cores=1  
--conf spark.executor.memory=4g --conf spark.driver.cores=1 --conf  
spark.driver.memory=4g --conf spark.executor.instances=1"  
    }  
  }'
```

3. Anote o ID de execução do trabalho retornado na saída. Substitua *job-run-id* por esse ID nas etapas a seguir.

Hive

Neste tutorial, criamos uma tabela, inserimos alguns registros e executamos uma consulta de agregação de contagem. Para executar o trabalho do Hive, primeiro crie um arquivo que contenha todas as consultas do Hive a serem executadas como parte de um único trabalho, faça upload do arquivo no S3 e especifique esse caminho do S3 ao iniciar o trabalho do Hive.

Para executar um trabalho do Hive

1. Crie um arquivo chamado `hive-query.q1` que contenha todas as consultas que você deseja executar no trabalho do Hive.

```
create database if not exists emrserverless;  
use emrserverless;  
create table if not exists test_table(id int);  
drop table if exists Values__Tmp__Table__1;  
insert into test_table values (1),(2),(2),(3),(3),(3);  
select id, count(id) from test_table group by id order by id desc;
```

2. Faça upload de `hive-query.q1` no bucket do S3 com o comando a seguir.

```
aws s3 cp hive-query.q1 s3://amzn-s3-demo-bucket/emr-serverless-hive/query/hive-  
query.q1
```

3. No comando a seguir, substitua *application-id* pelo ID da sua aplicação. Substitua *job-role-arn* pelo ARN do perfil de runtime que você criou em [Criação de um perfil de runtime de trabalhos](#). Substitua todas as strings *amzn-s3-demo-bucket* pelo bucket do Amazon S3 que você criou e adicione /output e /logs ao caminho. Isso cria pastas no bucket, nas quais o EMR Sem Servidor pode copiar os arquivos de saída e de log da aplicação.

```
aws emr-serverless start-job-run \  
  --application-id application-id \  
  --execution-role-arn job-role-arn \  
  --job-driver '{  
    "hive": {  
      "query": "s3://amzn-s3-demo-bucket/emr-serverless-hive/query/hive-  
query.q1",  
      "parameters": "--hiveconf hive.log.explain.output=false"  
    }  
  }' \  
  --configuration-overrides '{  
    "applicationConfiguration": [{  
      "classification": "hive-site",  
      "properties": {  
        "hive.exec.scratchdir": "s3://amzn-s3-demo-bucket/emr-serverless-  
hive/hive/scratch",  
        "hive.metastore.warehouse.dir": "s3://amzn-s3-demo-bucket/emr-  
serverless-hive/hive/warehouse",  
        "hive.driver.cores": "2",  
        "hive.driver.memory": "4g",  
        "hive.tez.container.size": "4096",  
        "hive.tez.cpu.vcores": "1"  
      }  
    }],  
    "monitoringConfiguration": {  
      "s3MonitoringConfiguration": {  
        "logUri": "s3://amzn-s3-demo-bucket/emr-serverless-hive/logs"  
      }  
    }  
  }'
```

4. Anote o ID de execução do trabalho retornado na saída. Substitua *job-run-id* por esse ID nas etapas a seguir.

Etapa 3: revisar a saída da execução do trabalho

A execução do trabalho normalmente leva de 3 a 5 minutos para ser concluída.

Spark

Você pode verificar o estado do trabalho do Spark com o comando a seguir.

```
aws emr-serverless get-job-run \  
  --application-id application-id \  
  --job-run-id job-run-id
```

Com o destino do log definido como `s3://amzn-s3-demo-bucket/emr-serverless-spark/logs`, você pode encontrar os logs desse trabalho específico executado em `s3://amzn-s3-demo-bucket/emr-serverless-spark/logs/applications/application-id/jobs/job-run-id`.

Em aplicações do Spark, o EMR Sem Servidor envia logs de eventos a cada 30 segundos para a pasta `sparklogs` no destino de log do S3. Quando o trabalho for concluído, os logs de runtime do Spark para o driver e os executores são enviados a pastas nomeadas apropriadamente pelo tipo de trabalhador, como `driver` ou `executor`. A saída do PySpark trabalho é enviada para o `s3://amzn-s3-demo-bucket/output/`

Hive

Você pode verificar o estado do trabalho do Hive com o comando a seguir.

```
aws emr-serverless get-job-run \  
  --application-id application-id \  
  --job-run-id job-run-id
```

Com o destino do log definido como `s3://amzn-s3-demo-bucket/emr-serverless-hive/logs`, você pode encontrar os logs desse trabalho específico executado em `s3://amzn-s3-demo-bucket/emr-serverless-hive/logs/applications/application-id/jobs/job-run-id`.

Em aplicações do Hive, o EMR Sem Servidor faz upload contínuo do driver do Hive para a pasta HIVE_DRIVER e os logs de tarefas do Tez para a pasta TEZ_TASK do destino de log do S3. Depois que a execução do trabalho atinge o estado SUCCEEDED, a saída da consulta do Hive fica disponível no local do Amazon S3 que você especificou no campo `monitoringConfiguration` de `configurationOverrides`.

Etapa 4: limpar

Após terminar de trabalhar com este tutorial, considere excluir os recursos criados por você. É recomendável liberar recursos que você não pretende usar novamente.

Exclua a aplicação

Para excluir uma aplicação, use o comando a seguir.

```
aws emr-serverless delete-application \  
  --application-id application-id
```

Exclusão do bucket de logs do S3

Para excluir o bucket de registro em log e saída do S3, use o comando a seguir. Substitua `amzn-s3-demo-bucket` pelo nome real do bucket do S3 criado em [Preparação do armazenamento do EMR Sem Servidor](#).

```
aws s3 rm s3://amzn-s3-demo-bucket --recursive  
aws s3api delete-bucket --bucket amzn-s3-demo-bucket
```

Exclusão do perfil de runtime do trabalho

Para excluir o perfil de runtime, separe a política do perfil. Em seguida, você pode excluir o perfil e a política.

```
aws iam detach-role-policy \  
  --role-name EMRServerlessS3RuntimeRole \  
  --policy-arn policy-arn
```

Para excluir o perfil, use o comando a seguir.

```
aws iam delete-role \  
  --role-name role-name
```

```
--role-name EMRServerlessS3RuntimeRole
```

Para excluir a política anexada ao perfil, use o comando a seguir.

```
aws iam delete-policy \  
  --policy-arn policy-arn
```

Para obter mais exemplos de execução de trabalhos do Spark e do Hive, consulte [Uso das configurações do Spark ao executar trabalhos do EMR Sem Servidor](#) e [Uso das configurações do Hive ao executar trabalhos do EMR Sem Servidor](#).

Como interagir e configurar uma aplicação do EMR Sem Servidor

Esta seção aborda como interagir com a aplicação do Amazon EMR Sem Servidor usando a AWS CLI. Ela também descreve a configuração de aplicações, a execução de personalizações e os padrões dos mecanismos Spark e Hive.

Tópicos

- [Estados da aplicação](#)
- [Criação de uma aplicação do EMR Sem Servidor usando o console do EMR Studio](#)
- [Interagindo com seu aplicativo EMR Serverless no AWS CLI](#)
- [Configuração de uma aplicação ao trabalhar com o EMR Sem Servidor](#)
- [Personalização de uma imagem do EMR Sem Servidor](#)
- [Configuração do acesso à VPC para que aplicações do EMR Sem Servidor se conectem aos dados](#)
- [Opções de arquitetura do Amazon EMR Sem Servidor](#)
- [Simultaneidade e enfileiramento de trabalhos para um aplicação do EMR Sem Servidor](#)

Estados da aplicação

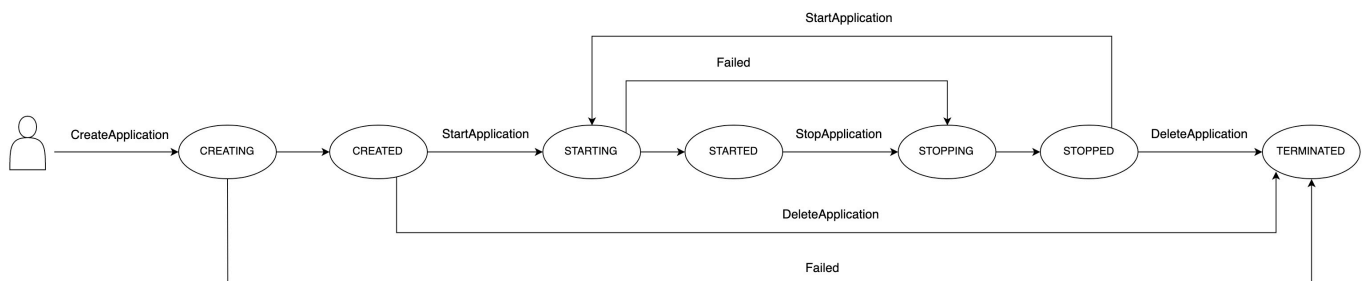
Quando você cria uma aplicação com o EMR Sem Servidor, a execução dela entra no estado CREATING. Em seguida, ela passa pelos estados apresentados a seguir até obter êxito (sair com o código 0) ou falhar (sair com um código diferente de zero).

As aplicações podem ter um dos seguintes status:

Estado	Description
Criando	A aplicação está sendo preparada e ainda não está pronta para uso.
Criado	A aplicação foi criada, mas ainda não provisionou a capacidade. Você pode modificar a

Estado	Description
	aplicação para alterar sua configuração de capacidade inicial.
Starting	A aplicação está sendo iniciada e provisionando a capacidade.
Iniciou	A aplicação está pronta para aceitar novos trabalhos. A aplicação só aceita trabalhos quando está nesse estado.
Stopping	Todos os trabalhos foram concluídos e a aplicação está liberando sua capacidade.
Interrompido	A aplicação foi interrompida e nenhum recurso está sendo executado na aplicação. Você pode modificar a aplicação para alterar sua configuração de capacidade inicial.
Encerrado	A aplicação foi encerrada e não aparece na sua lista de aplicações.

O diagrama a seguir mostra a trajetória dos estados da aplicação do EMR Sem Servidor.



Criação de uma aplicação do EMR Sem Servidor usando o console do EMR Studio

No console do EMR Studio, você pode criar, acessar e gerenciar aplicações do EMR Sem Servidor. Para navegar até o console do EMR Studio, siga as instruções em [Getting started from the console](#).

Criar uma aplicação do

Com a página Criar aplicação, crie uma aplicação do EMR Sem Servidor seguindo estas etapas.

1. No campo Nome, insira o nome que deseja chamar para a aplicação.
2. No campo Tipo, escolha Spark ou Hive como o tipo da aplicação.
3. No campo Versão de lançamento, escolha o número da versão do EMR.
4. Nas opções de Arquitetura, escolha a arquitetura do conjunto de instruções a ser usado. Para obter mais informações, consulte [Opções de arquitetura do Amazon EMR Sem Servidor](#).
 - arm64: arquitetura ARM de 64 bits para usar processadores Graviton
 - x86_64: arquitetura x86 de 64 bits para processadores baseados em x86
5. Há duas opções de configuração da aplicação para os campos restantes: configurações padrão e configurações personalizadas. Esses campos são opcionais.

Configurações padrão: as configurações padrão permitem criar rapidamente uma aplicação com capacidade pré-inicializada. Isso inclui um driver e um executor para o Spark e um driver e uma tarefa do Tez para o Hive. As configurações padrão não permitem a conectividade de rede com o seu VPCs. A aplicação está configurada para ser interrompida se ficar ociosa por 15 minutos e iniciar automaticamente no envio do trabalho.

Configurações personalizadas: as configurações personalizadas permitem modificar as propriedades a seguir.

- Capacidade pré-inicializada: o número de drivers e executores ou trabalhadores de tarefas do Tez do Hive e o tamanho de cada trabalhador.
- Limites da aplicação: a capacidade máxima de uma aplicação.
- Comportamento da aplicação: o comportamento de inicialização e interrupção automáticas da aplicação.
- Conexões de rede: conectividade de rede com recursos da VPC.

- Tags: tags personalizadas que atribui à aplicação.

Para obter mais informações sobre capacidade pré-inicializada, limites de aplicações e comportamento de aplicações, consulte [Configuração de uma aplicação ao trabalhar com o EMR Sem Servidor](#). Para obter mais informações sobre a conectividade de rede, consulte [Configuração do acesso à VPC para que aplicações do EMR Sem Servidor se conectem aos dados](#).

6. Para criar uma aplicação, escolha Criar aplicação.

Listagem de aplicações no console do EMR Studio

Você pode acessar todas as aplicações do EMR Sem Servidor existentes na página Lista de aplicações. Você pode escolher o nome de uma aplicação para navegar até a página Detalhes dela.

Gerenciamento de aplicações no console do EMR Studio

Você pode realizar as ações a seguir em uma aplicação na página Listagem de aplicações ou na página Detalhes de uma aplicação específica.

Iniciar a replicação

Escolha essa opção para iniciar manualmente uma aplicação.

Interromper a aplicação

Escolha essa opção para interromper manualmente uma aplicação. A aplicação não deve ter trabalhos em execução para ser interrompida. Para saber mais sobre as transições de estado de aplicações, consulte [Estados da aplicação](#).

Configuração de aplicações

Edite as configurações opcionais de uma aplicação na página Configuração de aplicações. Você pode alterar a maioria das configurações da aplicação. Por exemplo, altere o rótulo de lançamento de uma aplicação para atualizá-lo para uma versão diferente do Amazon EMR ou mude a arquitetura de x86_64 para arm64. As outras configurações opcionais são as mesmas que estão na seção Configurações personalizadas na página Criação de aplicações. Para obter mais informações sobre as configurações da aplicação, consulte [Criar uma aplicação do](#).

Exclusão de aplicação

Escolha essa opção para excluir manualmente uma aplicação. Você deve interromper uma aplicação para excluí-la. Para saber mais sobre as transições de estado de aplicações, consulte [Estados da aplicação](#).

Interagindo com seu aplicativo EMR Serverless no AWS CLI

A partir do AWS CLI, crie, descreva e exclua aplicativos individuais. Você também pode listar todos as aplicações a fim de acessá-las rapidamente. Esta seção descreve como executar essas ações. Para mais operações da aplicação, como iniciar, parar e atualizações da aplicação, consulte a [Referência da API do EMR Sem Servidor](#). Para obter exemplos de como usar a API do EMR Serverless usando o AWS SDK para Java, consulte [exemplos de Java](#) em nosso repositório. GitHub Para obter exemplos de como usar a API EMR Serverless usando o, AWS SDK para Python (Boto) consulte exemplos de [Python](#) em nosso repositório. GitHub

Para criar uma aplicação, use `create-application`. Você deve especificar SPARK ou HIVE como `type` de aplicação. Esse comando retorna o ARN, o nome e o ID da aplicação.

```
aws emr-serverless create-application \  
--name my-application-name \  
--type 'application-type' \  
--release-label release-version
```

Para descrever uma aplicação, use `get-application` e forneça o `application-id`. Esse comando retorna as configurações relacionadas ao estado e à capacidade da aplicação.

```
aws emr-serverless get-application \  
--application-id application-id
```

Para listas todas as suas aplicações, chame `list-applications`. Esse comando retorna as mesmas propriedades que `get-application`, mas inclui todas as aplicações.

```
aws emr-serverless list-applications
```

Para excluir a aplicação, chame `delete-application` e forneça o `application-id`.

```
aws emr-serverless delete-application \  

```

```
--application-id application-id
```

Configuração de uma aplicação ao trabalhar com o EMR Sem Servidor

Com o EMR Sem Servidor, configure as aplicações que você usa. Por exemplo, defina a capacidade máxima para a qual uma aplicação pode aumentar a escala verticalmente, configurar a capacidade pré-inicializada para manter o driver e os trabalhadores prontos para responder e especificar um conjunto comum de configurações de runtime e monitoramento no nível da aplicação. As páginas a seguir descrevem como configurar aplicações ao usar o EMR Sem Servidor.

Tópicos

- [Noções básicas do comportamento da aplicação no EMR Sem Servidor](#)
- [Capacidade pré-inicializada para trabalhar com uma aplicação no EMR Sem Servidor](#)
- [Configuração padrão de aplicações do EMR Sem Servidor](#)

Noções básicas do comportamento da aplicação no EMR Sem Servidor

Esta seção descreve o comportamento de envio de trabalhos, a configuração da capacidade para ajuste de escala e as definições de configuração do trabalhador do EMR Sem Servidor.

Comportamento padrão da aplicação

Início automático: por padrão, uma aplicação é configurada para iniciar automaticamente no envio do trabalho. Você pode desativar esse recurso.

Parada automática: por padrão, uma aplicação é configurada para ser interrompida automaticamente quando ociosa por 15 minutos. Quando uma aplicação muda para o estado STOPPED, ela libera qualquer capacidade pré-inicializada configurada. Você pode modificar a quantidade de tempo ocioso antes que uma aplicação pare automaticamente ou desativar esse recurso.

Capacidade máxima

Você pode configurar a capacidade máxima para a qual uma aplicação pode aumentar a escala verticalmente. Você pode especificar sua capacidade máxima em termos de CPU, memória (GB) e disco (GB).

Note

As práticas recomendadas são configurar sua capacidade máxima de forma proporcional aos tamanhos dos seus operadores com suporte, multiplicando o número de operadores pelos seus tamanhos. Por exemplo, se você quiser limitar seu aplicativo a 50 trabalhadores com 2 vCPUs, 16 GB para memória e 20 GB para disco, defina sua capacidade máxima para 100 vCPUs, 800 GB para memória e 1000 GB para disco.

Configuração de trabalhador compatíveis

A tabela a seguir mostra as configurações e tamanhos de trabalhadores compatíveis que você pode especificar para o EMR Sem Servidor. Configure tamanhos diferentes para drivers e executores com base na necessidade da workload.

Configurações e tamanhos de operadores

CPU	Memória	Armazenamento temporário padrão
1 vCPU	Mínimo de 2 GB, máximo de 8 GB, em incrementos de 1 GB	De 20 GB a 200 GB
2 vCPU	Mínimo de 4 GB, máximo de 16 GB, em incrementos de 1 GB	De 20 GB a 200 GB
4 vCPU	Mínimo de 8 GB, máximo de 30 GB, em incrementos de 1 GB	De 20 GB a 200 GB
8 vCPU	Mínimo de 16 GB, máximo de 60 GB, em incrementos de 4 GB	De 20 GB a 200 GB
16 vCPU	Mínimo de 32 GB, máximo de 120 GB, em incrementos de 8 GB	De 20 GB a 200 GB

CPU — Cada trabalhador pode ter 1, 2, 4, 8 ou 16 CPUs v.

Memória: cada trabalhador tem memória, especificada em GB, dentro dos limites listados na tabela anterior. Os trabalhos do Spark têm uma sobrecarga de memória, o que significa que a memória que eles usam é maior do que os tamanhos de contêineres especificados. Essa sobrecarga é especificada com as propriedades `spark.driver.memoryOverhead` e `spark.executor.memoryOverhead`. A sobrecarga tem um valor padrão de 10% da memória do contêiner, com um mínimo de 384 MB. Você deve considerar essa sobrecarga ao escolher o tamanho dos trabalhadores.

Por exemplo, se você escolher 4 v CPUs para sua instância de trabalho e uma capacidade de armazenamento pré-inicializada de 30 GB, defina um valor de aproximadamente 27 GB como memória executora para sua tarefa do Spark. Isso maximiza a utilização da capacidade pré-inicializada. A memória utilizável é 27 GB, mais 10% de 27 GB (2,7 GB), totalizando 29,7 GB.

Disco: você pode configurar cada trabalhador com discos de armazenamento temporário com tamanho mínimo de 20 GB e máximo de 200 GB. Você paga apenas pelo armazenamento adicional além de 20 GB configurado por trabalhador.

Capacidade pré-inicializada para trabalhar com uma aplicação no EMR Sem Servidor

O EMR Sem Servidor fornece um recurso opcional que mantém o driver e os trabalhadores pré-inicializados e prontos para responder em segundos. Isso cria efetivamente um grupo de aquecimento de trabalhadores para uma aplicação. Esse recurso é chamado de capacidade pré-inicializada. Para configurar esse recurso, defina o parâmetro `initialCapacity` de uma aplicação para o número de trabalhadores que você deseja pré-inicializar. Com a capacidade de trabalhadores pré-inicializada, os trabalhos começam imediatamente. Isso é ideal quando você deseja implementar aplicações iterativas e trabalhos urgentes.

A capacidade pré-inicializada mantém um grupo de aquecimento de trabalhadores prontos para a inicialização de trabalhos e sessões em segundos. Você pagará por trabalhadores pré-inicializados provisionados mesmo quando a aplicação estiver ociosa. Portanto, sugerimos habilitá-los para casos de uso que se beneficiem do rápido tempo de inicialização e dimensioná-los para a utilização ideal dos recursos. As aplicações do EMR Sem Servidor são desligadas automaticamente quando ociosas. Sugerimos manter esse recurso ativado ao usar trabalhadores pré-inicializados para evitar cobranças inesperadas.

Quando você envia um trabalho, se houver trabalhadores da `initialCapacity` disponíveis, o trabalho usa esses recursos para iniciar sua execução. Se esses trabalhadores já estiverem sendo usados por outros trabalhos, ou se o trabalho precisar de mais recursos do que os disponíveis da `initialCapacity`, a aplicação solicitará e obterá trabalhadores adicionais, até os limites máximos de recursos definidos para a aplicação. Quando um trabalho termina a execução, ele libera os trabalhadores que usou e o número de recursos disponíveis para a aplicação retorna para `initialCapacity`. Uma aplicação mantém a `initialCapacity` dos recursos mesmo após o término da execução dos trabalhos. A aplicação libera recursos em excesso além da `initialCapacity` quando os trabalhos não precisam mais deles para serem executados.

A capacidade pré-inicializada está disponível e pronta para uso quando a aplicação é iniciada. A capacidade pré-inicializada fica inativa quando a aplicação é interrompida. Uma aplicação passa para o estado `STARTED` somente se a capacidade pré-inicializada solicitada tiver sido criada e estiver pronta para uso. Durante todo o tempo em que a aplicação está no estado `STARTED`, o EMR Sem Servidor mantém a capacidade pré-inicializada disponível para utilização ou em uso por trabalhos ou workloads interativas. O recurso restaura a capacidade de contêineres liberados ou com falha. Isso mantém o número de trabalhadores que o parâmetro `InitialCapacity` especifica. O estado de uma aplicação sem capacidade pré-inicializada pode mudar imediatamente de `CREATED` para `STARTED`.

Você pode configurar a aplicação para liberar a capacidade pré-inicializada se ela não for usada por um determinado período, com um padrão de 15 minutos. Uma aplicação interrompida é iniciada automaticamente quando você envia um novo trabalho. Você pode definir essas configurações automáticas de início e parada ao criar a aplicação ou alterá-las quando a aplicação estiver em um estado `CREATED` ou `STOPPED`.

Você pode alterar as contagens da `InitialCapacity` e especificar configurações de computação, como CPU, memória e disco, para cada trabalhador. Como você não pode fazer modificações parciais, especifique todas as configurações de computação ao alterar os valores. Você só pode alterar as configurações quando a aplicação está no estado `CREATED` ou `STOPPED`.

Note

Para otimizar o uso de recursos da aplicação, sugerimos alinhar os tamanhos dos contêineres com os tamanhos dos trabalhadores da capacidade pré-inicializada. Por exemplo, se você configurar o tamanho do executor do Spark para 2 CPUs e sua memória para 8 GB, mas o tamanho do operador de capacidade pré-inicializada for 4 CPUs com

16 GB de memória, os executores do Spark usarão apenas metade dos recursos dos trabalhadores quando forem designados para esse trabalho.

Personalização da capacidade pré-inicializada para Spark e Hive

Você pode personalizar ainda mais a capacidade pré-inicializada de workloads executadas em estruturas específicas de big data. Por exemplo, quando uma workload é executada no Apache Spark, especifique quantos trabalhadores começam como drivers e quantos começam como executores. Da mesma forma, ao usar o Apache Hive, especifique quantos trabalhadores começam como drivers do Hive e quantos devem executar tarefas do Tez.

Configuração de aplicações executando o Apache Hive com capacidade pré-inicializada

A solicitação de API a seguir cria uma aplicação executando o Apache Hive com base na versão emr-6.6.0 do Amazon EMR. A aplicação começa com 5 drivers pré-inicializados do Hive, cada um com 2 vCPUs e 4 GB de memória, e 50 trabalhadores pré-inicializados de tarefas do Tez, cada um com 4 vCPUs e 8 GB de memória. Quando as consultas do Hive são executadas nessa aplicação, elas primeiro usam os trabalhadores pré-inicializados e começam a ser executadas imediatamente. Se todos os trabalhadores pré-inicializados estiverem ocupados e mais trabalhos do Hive forem enviados, a aplicação pode escalar para um total de 400 vCPU e 1.024 GB de memória. Opcionalmente, você pode omitir a capacidade do trabalhador do DRIVER ou da TEZ_TASK.

```
aws emr-serverless create-application \  
  --type "HIVE" \  
  --name my-application-name \  
  --release-label emr-6.6.0 \  
  --initial-capacity '{  
    "DRIVER": {  
      "workerCount": 5,  
      "workerConfiguration": {  
        "cpu": "2vCPU",  
        "memory": "4GB"  
      }  
    },  
    "TEZ_TASK": {  
      "workerCount": 50,  
      "workerConfiguration": {  
        "cpu": "4vCPU",  
        "memory": "8GB"  
      }  
    }  
  }'
```

```
    }  
  }  
}' \  
--maximum-capacity '{  
  "cpu": "400vCPU",  
  "memory": "1024GB"  
}'
```

Configuração de aplicações executando o Apache Spark com capacidade pré-inicializada

A solicitação de API a seguir cria uma aplicação que executa o Apache Spark 3.2.0 com base na versão 6.6.0 do Amazon EMR. A aplicação começa com 5 drivers pré-inicializados do Spark, cada um com 2 vCPUs e 4 GB de memória, e 50 executores pré-inicializados, cada um com 4 vCPUs e 8 GB de memória. Quando os trabalhos do Spark são executados nessa aplicação, eles primeiro usam os trabalhadores pré-inicializados e começam a ser executados imediatamente. Se todos os trabalhadores pré-inicializados estiverem ocupados e mais trabalhos do Spark forem enviados, a aplicação pode escalar para um total de 400 vCPU e 1.024 GB de memória. Opcionalmente, você pode omitir a capacidade do DRIVER ou do EXECUTOR.

Note

O Spark adiciona uma sobrecarga de memória configurável, com um valor padrão de 10%, à memória solicitada pelo driver e pelos executores. Para que os trabalhos usem trabalhadores pré-inicializados, a configuração inicial da memória de capacidade deve ser maior do que a memória solicitada pelo trabalho e pela sobrecarga.

```
aws emr-serverless create-application \  
--type "SPARK" \  
--name my-application-name \  
--release-label emr-6.6.0 \  
--initial-capacity '{  
  "DRIVER": {  
    "workerCount": 5,  
    "workerConfiguration": {  
      "cpu": "2vCPU",  
      "memory": "4GB"  
    }  
  },  
  "EXECUTOR": {
```

```
    "workerCount": 50,
    "workerConfiguration": {
      "cpu": "4vCPU",
      "memory": "8GB"
    }
  }
}' \
--maximum-capacity '{
  "cpu": "400vCPU",
  "memory": "1024GB"
}'
```

Configuração padrão de aplicações do EMR Sem Servidor

Você pode especificar um conjunto comum de configurações de runtime e monitoramento no nível da aplicação para todos os trabalhos enviados na mesma aplicação. Isso reduz a sobrecarga adicional associada à necessidade de enviar as mesmas configurações para cada trabalho.

Você pode modificar as configurações nos seguintes momentos:

- [Declare as configurações em nível de aplicação no envio do trabalho.](#)
- [Substitua as configurações padrão durante a execução do trabalho.](#)

As seções a seguir fornecem mais detalhes e um exemplo para contexto adicional.

Declaração de configurações no nível da aplicação

É possível especificar propriedades de configuração de registro em log e de runtime no nível da aplicação para os trabalhos enviados na aplicação.

monitoringConfiguration

Para especificar as configurações de log de trabalhos enviados com a aplicação, use o campo [monitoringConfiguration](#). Para obter mais informações sobre o registro em log do EMR Sem Servidor, consulte [Armazenamento de logs](#).

runtimeConfiguration

Para especificar propriedades de configuração de runtime, como `spark-defaults`, forneça um objeto de configuração no campo `runtimeConfiguration`. Isso afeta as configurações

padrão de todos os trabalhos enviados com a aplicação. Para obter mais informações, consulte [Parâmetro de substituição da configuração do Hive](#) e [Parâmetro de substituição de configuração do Spark](#).

As classificações de configuração disponíveis variam de acordo com a versão específica do EMR Sem Servidor. Por exemplo, classificações para Log4j personalizado `spark-driver-log4j2` e `spark-executor-log4j2` estão disponíveis somente nas versões 6.8.0 e superiores. Para obter uma lista de propriedades específicas da aplicação, consulte [Propriedades do trabalho do spark](#) e [Propriedades do trabalho do Hive](#).

Você também pode configurar [propriedades do Apache Log4j2](#), o [AWS Secrets Manager para proteção de dados](#) e o [runtime do Java 17](#) no nível da aplicação.

Para transmitir segredos do Secrets Manager no nível da aplicação, anexe a política a seguir aos usuários e perfis que precisam criar ou atualizar aplicações do EMR Sem Servidor com segredos.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "SecretsManagerPolicy",
      "Effect": "Allow",
      "Action": [
        "secretsmanager:GetSecretValue",
        "secretsmanager:DescribeSecret"
      ],
      "Resource": [
        "arn:aws:secretsmanager:us-east-1:123456789012:secret:my-secret-name-123abc"
      ]
    },
    {
      "Sid": "KMSDecryptPolicy",
      "Effect": "Allow",
      "Action": [
        "kms:Decrypt"
      ],
      "Resource": [
        "arn:aws:kms:us-east-1:123456789012:key/12345678-1234-1234-1234-123456789012"
      ]
    }
  ]
}
```

```

    ]
  }
]
}
```

Para obter mais informações sobre a criação de políticas personalizadas para segredos, consulte [Exemplos de políticas de permissões para AWS Secrets Manager](#) no Guia do usuário do AWS Secrets Manager .

Note

A `runtimeConfiguration` que você especifica no nível da aplicação é mapeada para `applicationConfiguration` na API [StartJobRun](#).

Exemplo de declaração

O exemplo a seguir mostra como declarar configurações padrão com `create-application`.

```
aws emr-serverless create-application \
  --release-label release-version \
  --type SPARK \
  --name my-application-name \
  --runtime-configuration '[
    {
      "classification": "spark-defaults",
      "properties": {
        "spark.driver.cores": "4",
        "spark.executor.cores": "2",
        "spark.driver.memory": "8G",
        "spark.executor.memory": "8G",
        "spark.executor.instances": "2",

        "spark.hadoop.java.x.jdo.option.ConnectionDriverName": "org.mariadb.jdbc.Driver",
        "spark.hadoop.java.x.jdo.option.ConnectionURL": "jdbc:mysql://db-host:db-
port/db-name",
        "spark.hadoop.java.x.jdo.option.ConnectionUserName": "connection-user-
name",
        "spark.hadoop.java.x.jdo.option.ConnectionPassword":
        "EMR.secret@SecretID"
      }
    }
  ]'
```

```

    },
    {
      "classification": "spark-driver-log4j2",
      "properties": {
        "rootLogger.level": "error",
        "logger.IdentifierForClass.name": "classpathForSettingLogger",
        "logger.IdentifierForClass.level": "info"
      }
    }
  ]' \
  --monitoring-configuration '{
    "s3MonitoringConfiguration": {
      "logUri": "s3://amzn-s3-demo-logging-bucket/logs/app-level"
    },
    "managedPersistenceMonitoringConfiguration": {
      "enabled": false
    }
  }'

```

Substituição de configurações durante a execução de um trabalho

Você pode especificar substituições de configuração para a configuração da aplicação e a configuração de monitoramento com a API [StartJobRun](#). Em seguida, o EMR Sem Servidor mescla as configurações que você especifica no nível da aplicação e no nível do trabalho para determinar as configurações da execução do trabalho.

O nível de granularidade quando a mesclagem ocorre é o seguinte:

- **[ApplicationConfiguration](#)**: tipo de classificação, por exemplo, spark-defaults.
- **[MonitoringConfiguration](#)**: tipo de configuração, por exemplo, s3MonitoringConfiguration.

Note

A prioridade das configurações fornecidas em [StartJobRun](#) substitui as configurações que você fornece no nível da aplicação.

Para obter mais informações sobre classificações de prioridade, consulte [Parâmetro de substituição da configuração do Hive](#) e [Parâmetro de substituição de configuração do Spark](#).

Quando você inicia um trabalho, caso não especifique uma configuração específica, ela será herdada da aplicação. Se você declarar as configurações no nível do trabalho, poderá executar as seguintes operações:

- Substituir uma configuração existente: forneça o mesmo parâmetro de configuração na solicitação de `StartJobRun` com seus valores de substituição.
- Adicionar mais uma configuração: adicione o novo parâmetro de configuração na solicitação de `StartJobRun` com os valores que deseja especificar.
- Remover uma configuração existente: para remover a configuração de runtime de uma aplicação, forneça a chave da configuração que deseja remover e passe uma declaração vazia `{}` para a configuração. Não recomendamos remover nenhuma classificação que contenha parâmetros necessários para a execução de um trabalho. Por exemplo, se você tentar remover as [propriedades necessárias de um trabalho do Hive](#), o trabalho falhará.

Para remover a configuração de monitoramento de uma aplicação, use o método apropriado para o tipo de configuração relevante:

- **cloudWatchLoggingConfiguration**: para remover `cloudWatchLogging`, passe o sinalizador habilitado como `false`.
- **managedPersistenceMonitoringConfiguration**: para remover as configurações de persistência gerenciada e voltar ao estado habilitado padrão, passe uma declaração vazia `{}` para a configuração.
- **s3MonitoringConfiguration**: para remover `s3MonitoringConfiguration`, passe uma declaração vazia `{}` para a configuração.

Exemplo de substituição

O exemplo a seguir mostra diferentes operações que você pode realizar durante o envio do trabalho em `start-job-run`.

```
aws emr-serverless start-job-run \  
  --application-id your-application-id \  
  --execution-role-arn your-job-role-arn \  
  --job-driver '{  
    "sparkSubmit": {  
      "entryPoint": "s3://us-east-1.elasticmapreduce/emr-containers/samples/  
wordcount/scripts/wordcount.py",  
      "entryPointArguments": ["s3://amzn-s3-demo-destination-bucket1/  
wordcount_output"]
```

```

    }
  }' \
  --configuration-overrides '{
    "applicationConfiguration": [
      {
        // Override existing configuration for spark-defaults in the
application
        "classification": "spark-defaults",
        "properties": {
          "spark.driver.cores": "2",
          "spark.executor.cores": "1",
          "spark.driver.memory": "4G",
          "spark.executor.memory": "4G"
        }
      },
      {
        // Add configuration for spark-executor-log4j2
        "classification": "spark-executor-log4j2",
        "properties": {
          "rootLogger.level": "error",
          "logger.IdentifierForClass.name": "classpathForSettingLogger",
          "logger.IdentifierForClass.level": "info"
        }
      },
      {
        // Remove existing configuration for spark-driver-log4j2 from the
application
        "classification": "spark-driver-log4j2",
        "properties": {}
      }
    ],
    "monitoringConfiguration": {
      "managedPersistenceMonitoringConfiguration": {
        // Override existing configuration for managed persistence
        "enabled": true
      },
      "s3MonitoringConfiguration": {
        // Remove configuration of S3 monitoring
      },
      "cloudWatchLoggingConfiguration": {
        // Add configuration for CloudWatch logging
        "enabled": true
      }
    }
  }

```

```
}'
```

No momento da execução do trabalho, as classificações e configurações a seguir serão aplicadas com base na classificação de substituição de prioridade descrita em [Parâmetro de substituição da configuração do Hive](#) e [Parâmetro de substituição de configuração do Spark](#).

- A classificação `spark-defaults` será atualizada com as propriedades especificadas no nível do trabalho. Somente as propriedades incluídas em `StartJobRun` são consideradas para essa classificação.
- A classificação `spark-executor-log4j2` será adicionada à lista existente de classificações.
- A classificação `spark-driver-log4j2` será removida.
- As configurações de `managedPersistenceMonitoringConfiguration` serão atualizadas com as configurações no nível do trabalho.
- As configurações de `s3MonitoringConfiguration` serão removidas.
- As configurações de `cloudWatchLoggingConfiguration` serão adicionadas às configurações de monitoramento existentes.

Personalização de uma imagem do EMR Sem Servidor

A partir do Amazon EMR 6.9.0, use imagens personalizadas para empacotar dependências de aplicações e ambientes de runtime em um único contêiner com o Amazon EMR Sem Servidor. Isso simplifica a forma como você gerencia as dependências da workload e torna seus pacotes mais portáteis. Quando você personaliza a imagem do EMR Sem Servidor, ela oferece os seguintes benefícios:

- Instala e configura pacotes otimizados para as workloads. Esses pacotes não estão amplamente disponíveis na distribuição pública dos ambientes de runtime do Amazon EMR.
- Integra o EMR Sem Servidor aos processos atuais de criação, teste e implantação estabelecidos em sua organização, incluindo desenvolvimento e testes locais.
- Aplica os processos de segurança estabelecidos, como a verificação de imagens, que atendem aos requisitos de conformidade e governança da sua organização.
- Permite que você use suas próprias versões do JDK e do Python em aplicações.

O EMR Sem Servidor fornece imagens que você usa como base ao criar suas próprias imagens. A imagem base fornece os jars, a configuração e as bibliotecas essenciais para que a imagem interaja

com o EMR Sem Servidor. É possível encontrar a imagem base na [Galeria pública do Amazon ECR](#). Use a imagem que corresponda ao seu tipo de aplicação (Spark ou Hive) e à versão de lançamento. Por exemplo, se você criar uma aplicação no Amazon EMR 6.9.0, use as imagens a seguir.

Tipo	Imagem
Spark	public.ecr.aws/emr-serverless/spark/emr-6.9.0:latest
Hive	public.ecr.aws/emr-serverless/hive/emr-6.9.0:latest

Pré-requisitos

Antes de criar uma imagem personalizada do EMR Sem Servidor, preencha estes pré-requisitos.

1. Crie um repositório Amazon ECR no mesmo Região da AWS que você usa para iniciar aplicativos EMR Serverless. Para criar um repositório privado do Amazon ECR, consulte [Criar um repositório privado](#).
2. Para conceder aos usuários acesso ao repositório do Amazon ECR, adicione as políticas a seguir aos usuários e perfis que criam ou atualizam aplicações do EMR Sem Servidor com imagens desse repositório.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "ECRRepositoryListGetPolicy",
      "Effect": "Allow",
      "Action": [
        "ecr:GetDownloadUrlForLayer",
        "ecr:BatchGetImage",
        "ecr:DescribeImages"
      ],
      "Resource": [
        "arn:aws:ecr:*:123456789012:repository/my-repo"
      ]
    }
  ]
}
```

```
]
}
```

Para obter mais exemplos de políticas baseadas em identidade do Amazon ECR, consulte [Exemplos de políticas baseadas em identidade do Amazon Elastic Container Registry](#).

Etapa 1: criar uma imagem personalizada usando imagens base do EMR Sem Servidor

Primeiro, crie um [Dockerfile](#) que comece com uma instrução FROM que usa sua imagem base preferencial. Após a instrução FROM, inclua qualquer modificação que desejar fazer na imagem. A imagem base define automaticamente o valor de USER para hadoop. Essa configuração não tem permissões para todas as modificações que você inclui. Como solução, defina o USER para root, modifique sua imagem e, em seguida, defina o USER de volta para hadoop: hadoop. Para consultar exemplos de casos de uso comuns, veja [Uso de imagens personalizadas com o EMR Sem Servidor](#).

```
# Dockerfile
FROM public.ecr.aws/emr-serverless/spark/emr-6.9.0:latest

USER root
# MODIFICATIONS GO HERE

# EMRS runs the image as hadoop
USER hadoop:hadoop
```

Depois que você tiver o Dockerfile, compile a imagem com o comando a seguir.

```
# build the docker image
docker build . -t aws-account-id.dkr.ecr.region.amazonaws.com/my-repository[:tag]or[@digest]
```

Etapa 2: validar a imagem localmente

O EMR Sem Servidor fornece uma ferramenta off-line que pode verificar estaticamente sua imagem personalizada para validar arquivos base, variáveis de ambiente e configurações corretas de imagem. Para obter informações sobre como instalar e executar a ferramenta, consulte a CLI de imagem [sem servidor do Amazon EMR. GitHub](#)

Depois de instalar a ferramenta, execute o comando a seguir para validar uma imagem:

```
amazon-emr-serverless-image \
validate-image -r emr-6.9.0 -t spark \
-i aws-account-id.dkr.ecr.region.amazonaws.com/my-repository:tag/@digest
```

A saída é semelhante à seguinte.

```
Amazon EMR Serverless - Image CLI
Version: 0.0.1
... Checking if docker cli is installed
... Checking Image Manifest
[INFO] Image ID: 9e2f4359cf5beb466a8a2ed047ab61c9d37786c555655fc122272758f761b41a
[INFO] Created On: 2022-12-02T07:46:42.586249984Z
[INFO] Default User Set to hadoop:hadoop : PASS
[INFO] Working Directory Set to : PASS
[INFO] Entrypoint Set to /usr/bin/entrypoint.sh : PASS
[INFO] HADOOP_HOME is set with value: /usr/lib/hadoop : PASS
[INFO] HADOOP_LIBEXEC_DIR is set with value: /usr/lib/hadoop/libexec : PASS
[INFO] HADOOP_USER_HOME is set with value: /home/hadoop : PASS
[INFO] HADOOP_YARN_HOME is set with value: /usr/lib/hadoop-yarn : PASS
[INFO] HIVE_HOME is set with value: /usr/lib/hive : PASS
[INFO] JAVA_HOME is set with value: /etc/alternatives/jre : PASS
[INFO] TEZ_HOME is set with value: /usr/lib/tez : PASS
[INFO] YARN_HOME is set with value: /usr/lib/hadoop-yarn : PASS
[INFO] File Structure Test for hadoop-files in /usr/lib/hadoop: PASS
[INFO] File Structure Test for hadoop-jars in /usr/lib/hadoop/lib: PASS
[INFO] File Structure Test for hadoop-yarn-jars in /usr/lib/hadoop-yarn: PASS
[INFO] File Structure Test for hive-bin-files in /usr/bin: PASS
[INFO] File Structure Test for hive-jars in /usr/lib/hive/lib: PASS
[INFO] File Structure Test for java-bin in /etc/alternatives/jre/bin: PASS
[INFO] File Structure Test for tez-jars in /usr/lib/tez: PASS
-----
Overall Custom Image Validation Succeeded.
-----
```

Etapa 3: upload da imagem em um repositório do Amazon ECR

Envie a imagem do Amazon ECR para o repositório do Amazon ECR com os comandos a seguir. Verifique se você tem as permissões corretas do IAM para enviar a imagem ao repositório. Para obter mais informações, consulte [Enviar uma imagem](#) no Guia do usuário do Amazon ECR.

```
# login to ECR repo
```

```
aws ecr get-login-password --region region | docker login --username AWS --password-stdin aws-account-id.dkr.ecr.region.amazonaws.com

# push the docker image
docker push aws-account-id.dkr.ecr.region.amazonaws.com/my-repository:tag/@digest
```

Etapa 4: criar ou atualizar uma aplicação com imagens personalizadas

Escolha a Console de gerenciamento da AWS guia ou AWS CLI guia de acordo com a forma como você deseja iniciar seu aplicativo e, em seguida, conclua as etapas a seguir.

Console

1. [Faça login no console do EMR Studio em `https://console.aws.amazon.com/emr`](https://console.aws.amazon.com/emr). Navegue até a aplicação ou crie uma com as instruções em [Criação de aplicações](#).
2. Para especificar imagens personalizadas ao criar ou atualizar uma aplicação do EMR Sem Servidor, selecione Configurações personalizadas nas opções de configuração da aplicação.
3. Na seção Configurações de imagem personalizada, marque a caixa de seleção Usar a imagem personalizada com esta aplicação.
4. Cole o URI da imagem do Amazon ECR no campo URI da imagem. O EMR Sem Servidor usa essa imagem para todos os tipos de trabalhadores da aplicação. Como alternativa, você pode escolher imagens personalizadas diferentes e colar imagens diferentes do Amazon ECR URIs para cada tipo de trabalhador.

CLI

- Crie uma aplicação com o parâmetro `image-configuration`. O EMR Sem Servidor aplica essa configuração a todos os tipos de trabalhadores.

```
aws emr-serverless create-application \
--release-label emr-6.9.0 \
--type SPARK \
--image-configuration '{
  "imageUri": "aws-account-id.dkr.ecr.region.amazonaws.com/my-repository:tag/@digest"
}'
```

Para criar uma aplicação com configurações de imagem diferentes para cada tipo de trabalhador, use o parâmetro `worker-type-specifications`.

```
aws emr-serverless create-application \
--release-label emr-6.9.0 \
--type SPARK \
--worker-type-specifications '{
  "Driver": {
    "imageConfiguration": {
      "imageUri": "aws-account-id.dkr.ecr.region.amazonaws.com/my-
repository:tag/@digest"
    }
  },
  "Executor" : {
    "imageConfiguration": {
      "imageUri": "aws-account-id.dkr.ecr.region.amazonaws.com/my-
repository:tag/@digest"
    }
  }
}'
```

Para atualizar uma aplicação, use o parâmetro `image-configuration`. O EMR Sem Servidor aplica essa configuração a todos os tipos de trabalhadores.

```
aws emr-serverless update-application \
--application-id application-id \
--image-configuration '{
  "imageUri": "aws-account-id.dkr.ecr.region.amazonaws.com/my-repository:tag/
@digest"
}'
```

Etapa 5: permitir que o EMR Sem Servidor acesse o repositório de imagens personalizadas

Adicione a política de recursos a seguir ao repositório do Amazon ECR para permitir que a entidade principal do serviço do EMR Sem Servidor use as solicitações `get`, `describe` e `download` desse repositório.

JSON

```
{
```

```

"Version": "2012-10-17",
"Statement": [
  {
    "Sid": "EmrServerlessCustomImageSupport",
    "Effect": "Allow",
    "Action": [
      "ecr:BatchGetImage",
      "ecr:DescribeImages",
      "ecr:GetDownloadUrlForLayer"
    ],
    "Resource": "arn:aws:ecr:*:123456789012:repository/my-repo",
    "Condition": {
      "ArnLike": {
        "aws:SourceArn": "arn:aws:emr-serverless:*:123456789012:/applications/
*"
      }
    }
  }
]
}

```

Como uma prática recomendada de segurança, adicione uma chave de condição `aws:SourceArn` à política do repositório. A chave de condição global do IAM `aws:SourceArn` garante que o EMR Sem Servidor use o repositório somente para o ARN de uma aplicação. Para obter mais informações sobre as políticas de repositório do Amazon ECR, consulte [Criando um repositório privado](#).

Considerações e limitações

Considere o seguinte ao trabalhar com imagens personalizadas:

- Use a imagem base correta que corresponda ao tipo (Spark ou Hive) e ao rótulo de lançamento (por exemplo, `emr-6.9.0`) da aplicação.
- O EMR Sem Servidor ignora as instruções `[CMD]` ou `[ENTRYPOINT]` no arquivo Docker. Use instruções comuns no arquivo Docker, como `[COPY]`, `[RUN]` e `[WORKDIR]`.
- Não modifique as variáveis de ambiente `JAVA_HOME`, `SPARK_HOME`, `HIVE_HOME` e `TEZ_HOME` ao criar uma imagem personalizada.
- As imagens personalizadas não podem exceder 10 GB de tamanho.
- Se você modificar binários ou jars nas imagens base do Amazon EMR, isso poderá causar falhas no lançamento de aplicações ou trabalhos.

- O repositório do Amazon ECR deve estar no mesmo Região da AWS que você usa para iniciar aplicativos EMR Serverless.

Configuração do acesso à VPC para que aplicações do EMR Sem Servidor se conectem aos dados

Você pode configurar aplicações do EMR Sem Servidor para se conectar aos seus armazenamentos de dados na VPC, como clusters do Amazon Redshift, bancos de dados do Amazon RDS ou buckets do Amazon S3 com endpoints da VPC. A aplicação do EMR Sem Servidor tem conectividade de saída com os armazenamentos de dados na VPC. Por padrão, o EMR Sem Servidor bloqueia tanto o acesso de entrada às suas aplicações quanto o acesso de saída à Internet para reforçar a segurança.

Note

Você deve configurar o acesso à VPC se quiser usar um banco de dados externo do Hive Metastore para a aplicação. Para obter informações sobre como configurar um metastore externa do Hive, consulte [Configuração do metastore](#).

Criar aplicação

Na página Criar aplicação, escolha configurações personalizadas e especificar a VPC, as sub-redes e os grupos de segurança que as aplicações do EMR Sem Servidor podem usar.

VPCs


Escolha o nome da nuvem privada virtual (VPC) que contém os armazenamentos de dados. A página Criar aplicativo lista todas as VPCs opções escolhidas Região da AWS.

Sub-redes

Escolha as sub-redes na VPC que contém seu armazenamento de dados. A página Criação de aplicações lista todas as sub-redes dos armazenamentos de dados na VPC. Sub-redes públicas e privadas têm suporte. Você pode passar sub-redes privadas ou públicas para suas aplicações. A escolha entre uma sub-rede pública ou privada envolve algumas considerações importantes que devem ser levadas em conta.

Para sub-redes privadas:

- As tabelas de rotas associadas não devem ter gateways da Internet.
- Para conectividade de saída para a Internet, se necessário, configure rotas de saída usando um Gateway de NAT. Para configurar um gateway de NAT, consulte [Gateways de NAT](#).
- Para conectividade com o Amazon S3, configure um Gateway de NAT ou endpoint de VPC. Para configurar um endpoint da VPC do S3, consulte [Criar um endpoint do gateway](#).
- Se você configurar um endpoint de VPC do S3 e anexar uma política de endpoint para controlar o acesso, siga as instruções em [Registro em log para o EMR Sem Servidor com armazenamento gerenciado](#) para fornecer permissões para o EMR Sem Servidor armazenar e servir logs de aplicações.
- Para conectividade com outras pessoas Serviços da AWS fora da VPC, como com o Amazon DynamoDB, configure os VPC endpoints ou um gateway NAT. Para configurar endpoints da VPC de Serviços da AWS, consulte [Trabalhar com endpoints de VPC](#).

 Note

Ao configurar uma aplicação do Amazon EMR Sem Servidor em uma sub-rede privada, sugerimos que você também configure endpoints de VPC para o Amazon S3. Se a sua aplicação do EMR Sem Servidor estiver em uma sub-rede privada sem endpoints de VPC para o Amazon S3, haverá cobranças adicionais do gateway de NAT associadas ao tráfego do S3. Isso acontece porque o tráfego entre sua aplicação do EMR e o Amazon S3 não permanecerá dentro da sua VPC quando os endpoints de VPC não estiverem configurados.

Para sub-redes públicas:

- Eles têm uma rota para um Gateway da Internet.
- Você deve garantir configurações adequadas do grupo de segurança para controlar o tráfego de saída.

Os trabalhadores podem se conectar aos armazenamentos de dados na sua VPC por meio do tráfego de saída. Por padrão, o EMR Sem Servidor bloqueia o acesso de entrada aos operadores. Isso é para reforçar a segurança.

Quando você usa AWS Config, o EMR Serverless cria um registro de item de interface de rede elástica para cada trabalhador. Para evitar custos relacionados a esse recurso, considere desligá-lo `AWS::EC2::NetworkInterface` AWS Config.

Note

Sugerimos selecionar várias sub-redes entre várias zonas de disponibilidade. Isso ocorre porque as sub-redes que você escolhe determinam as zonas de disponibilidade que estão disponíveis para a inicialização de uma aplicação do EMR Sem Servidor. Cada trabalhador consome um endereço IP na sub-rede em que é iniciado. Certifique-se de que as sub-redes especificadas tenham endereços IP suficientes para o número de trabalhadores que você planeja iniciar. Para obter mais informações sobre o planejamento de sub-redes, consulte [the section called “Práticas recomendadas do planejamento de sub-redes”](#).

Considerações e limitações para sub-redes

- O EMR Serverless com sub-redes públicas não é compatível com Lake Formation. AWS
- O tráfego de entrada não é compatível com sub-redes públicas.

Grupos de segurança

Escolha um ou mais grupos de segurança que possam se comunicar com seus armazenamentos de dados. A página Criação de aplicações lista todos os grupos de segurança na sua VPC. O EMR Sem Servidor associa esses grupos de segurança a interfaces de rede elástica anexadas às sub-redes da VPC.

Note

Sugerimos criar um grupo de segurança separado para aplicações do EMR Sem Servidor. O EMR Sem Servidor não permite que você crie/atualize/inicie a aplicação se os grupos de segurança tiverem portas abertas para a Internet pública em 0.0.0.0/0 ou no intervalo ::/0. Isso proporciona maior segurança, isolamento e torna o gerenciamento das regras de rede mais eficiente. Por exemplo, isso bloqueia o tráfego inesperado para trabalhadores com endereços IP públicos. Para se comunicar com clusters do Amazon Redshift, por exemplo, defina as regras de tráfego entre os grupos de segurança do Redshift e do EMR Sem Servidor, conforme demonstrado no exemplo na seção a seguir.

Example Exemplo: comunicação com clusters do Amazon Redshift

1. Adicione uma regra para tráfego de entrada ao grupo de segurança do Amazon Redshift em um dos grupos de segurança do EMR Sem Servidor.

Tipo	Protocolo	Intervalo de portas	Fonte
Todos os TCP	TCP	5439	emr-serve rless-sec urity-group

2. Adicione uma regra para o tráfego de saída em um dos grupos de segurança do EMR Sem Servidor. É possível fazer isso de duas maneiras. Primeiro, abra o tráfego de saída para todas as portas.

Tipo	Protocolo	Intervalo de portas	Destino
Todo o tráfego	TCP	ALL	0.0.0.0/0

Como alternativa, você pode restringir o tráfego de saída para os clusters do Amazon Redshift. Isso é útil somente quando a aplicação precisa se comunicar com os clusters do Amazon Redshift e nada mais.

Tipo	Protocolo	Intervalo de portas	Fonte
Todos os TCP	TCP	5439	redshift- security- group

Configuração de aplicações

Você pode alterar a configuração da rede de uma aplicação do EMR Sem Servidor existente na página Configuração de aplicações.

Exibição dos detalhes da execução do trabalho

Na página de Detalhes da execução do trabalho, acesse a sub-rede usada pelo seu trabalho para uma execução específica. Observe que um trabalho é executado somente em uma sub-rede selecionada das sub-redes especificadas.

Práticas recomendadas do planejamento de sub-redes

AWS os recursos são criados em uma sub-rede que é um subconjunto de endereços IP disponíveis em uma Amazon VPC. Por exemplo, uma VPC com uma máscara de rede /16 tem até 65.536 endereços IP disponíveis que podem ser divididos em várias redes menores usando máscaras de sub-rede. Como exemplo, você pode dividir esse intervalo em duas sub-redes, cada uma usando a máscara /17 e 32.768 endereços IP disponíveis. Uma sub-rede reside dentro de uma zona de disponibilidade e não pode abranger várias zonas.

As sub-redes devem ser projetadas tendo em mente os limites de ajuste de escala da aplicação do EMR Sem Servidor. Por exemplo, se você tiver uma aplicação solicitando 4 trabalhadores de vCPU e puder aumentar a escala verticalmente para até 4.000 vCPUs, sua aplicação precisa de no máximo 1.000 trabalhadores para um total de 1.000 interfaces de rede. Sugerimos criar sub-redes em várias zonas de disponibilidade. Isso permite que o EMR Sem Servidor tente novamente seu trabalho ou provisione a capacidade pré-inicializada em uma zona de disponibilidade diferente em um evento improvável quando uma zona de disponibilidade falhar. Portanto, cada sub-rede em pelo menos duas zonas de disponibilidade deve ter mais de 1.000 endereços IP disponíveis.

Você precisa de sub-redes com tamanho de máscara menor ou igual a 22 para provisionar 1.000 interfaces de rede. Qualquer máscara maior que 22 não atende ao requisito. Por exemplo, uma máscara de sub-rede de /23 fornece 512 endereços IP, enquanto uma máscara de /22 fornece 1.024 e uma máscara de /21 fornece 2.048 endereços IP. Abaixo está um exemplo de 4 sub-redes com máscara de /22 em uma VPC de máscara de rede de /16 que podem ser alocadas para diferentes zonas de disponibilidade. Há uma diferença de cinco entre endereços IP disponíveis e utilizáveis porque os primeiros quatro endereços IP e o último endereço IP em cada sub-rede são reservados por. AWS

ID da sub-rede	Endereço de sub-rede	Máscara de sub-rede	Intervalo de endereços IP	Endereços IP disponíveis	Endereços IP utilizáveis
1	10.0.0.0	255.255.252.0/22	De 10.0.0.0 a 10.0.3.255	1,024	1.019

ID da sub-rede	Endereço de sub-rede	Máscara de sub-rede	Intervalo de endereços IP	Endereços IP disponíveis	Endereços IP utilizáveis
2	10.0.4.0	255.255.252.0/22	De 10.0.4.0 a 10.0.7.255	1,024	1.019
3	10.0.8.0	255.255.252.0/22	10.0.8.0 - 10.0.11.255	1,024	1.019
4	10.0.12.0	255.255.252.0/22	De 10.0.12.0 a 10.0.15.255	1,024	1.019

Você deve avaliar se a workload é mais adequada para trabalhadores maiores. Usar trabalhadores maiores requer menos interfaces de rede. Por exemplo, usar trabalhadores de 16 vCPUs com um limite de ajuste de escala de aplicações de 4.000 vCPUs exige no máximo 250 trabalhadores para um total de 250 endereços IP disponíveis para provisionar interfaces de rede. Você precisa de sub-redes em várias zonas de disponibilidade com tamanho de máscara menor ou igual a 24 para provisionar 250 interfaces de rede. Qualquer tamanho de máscara maior que 24 oferece menos de 250 endereços IP.

Se você compartilha sub-redes em várias aplicações, cada sub-rede deve ser projetada tendo em mente os limites de ajuste de escala coletivo de todas as aplicações. Por exemplo, se você tiver 3 aplicações solicitando 4 trabalhadores de vCPU e cada um puder aumentar a escala verticalmente para até 4.000 vCPUs com uma cota baseada em serviço em nível de conta de 12.000 vCPUs, cada sub-rede exige 3.000 endereços IP disponíveis. Se a VPC que você deseja utilizar não tiver uma quantidade suficiente de endereços IP, tente aumentar a quantidade de endereços IP disponíveis. É possível fazer isso associando blocos de Encaminhamento Entre Domínios Sem Classificação (CIDR) com sua VPC. Para obter mais informações, consulte [Associar blocos IPv4 CIDR adicionais à sua VPC no Guia](#) do usuário da Amazon VPC.

Você pode usar uma das muitas ferramentas disponíveis on-line para gerar rapidamente definições de sub-rede e analisar o intervalo disponível de endereços IP.

Opções de arquitetura do Amazon EMR Sem Servidor

A arquitetura do conjunto de instruções da aplicação do Amazon EMR Sem Servidor determina o tipo de processador que a aplicação usa para executar o trabalho. O Amazon EMR fornece

duas opções de arquitetura para a aplicação: x86_64 e arm64. O EMR Sem Servidor é atualizado automaticamente para a última geração de instâncias à medida que elas se tornam disponíveis, para que suas aplicações possam usar as instâncias mais novas sem exigir esforço adicional de sua parte.

Tópicos

- [Uso da arquitetura x86_64](#)
- [Uso da arquitetura arm64 \(Graviton\)](#)
- [Lançamento de aplicações com suporte ao Graviton](#)
- [Configuração de aplicações existentes para usar o Graviton](#)
- [Considerações ao usar o Graviton](#)

Uso da arquitetura x86_64

A arquitetura x86_64 também é conhecida como x86 de 64 bits ou x64. x86_64 é a opção padrão para aplicações do EMR Sem Servidor. Essa arquitetura usa processadores baseados em x86 e é compatível com a maioria das ferramentas e bibliotecas de terceiros.

A maioria das aplicações é compatível com a plataforma de hardware x86 e pode ser executada com êxito na arquitetura x86_64 padrão. No entanto, se a aplicação for compatível com ARM de 64 bits, mude para arm64 e usar os processadores Graviton para melhorar a performance, a potência computacional e a memória. Custa menos executar instâncias na arquitetura arm64 do que executar instâncias do mesmo tamanho na arquitetura x86.

Uso da arquitetura arm64 (Graviton)

AWS Os processadores Graviton são projetados de forma personalizada AWS com núcleos ARM Neoverse de 64 bits e utilizam a arquitetura arm64 (também conhecida como Arch64 ou ARM de 64 bits). A linha de processadores AWS Graviton disponível no EMR Serverless inclui os processadores Graviton3 e Graviton2. Esses processadores oferecem um melhor custo-benefício para workloads do Spark e do Hive em comparação com workloads equivalentes executadas na arquitetura x86_64. O EMR Sem Servidor usa automaticamente a última geração de processadores quando disponível, sem nenhum esforço de sua parte para fazer o upgrade para a última geração de processadores.

Lançamento de aplicações com suporte ao Graviton

Use um dos métodos a seguir para iniciar uma aplicação com a arquitetura arm64.

AWS CLI

Para iniciar um aplicativo usando processadores Graviton a partir de AWS CLI, especifique ARM64 como `architecture` parâmetro na `create-application` API. Forneça os valores adequados para a aplicação em outros parâmetros.

```
aws emr-serverless create-application \  
  --name my-graviton-app \  
  --release-label emr-6.8.0 \  
  --type "SPARK" \  
  --architecture "ARM64" \  
  --region us-west-2
```

EMR Studio

Para iniciar uma aplicação usando os processadores Graviton do EMR Studio, escolha `arm64` como a opção `Arquitetura` ao criar ou atualizar uma aplicação.

Configuração de aplicações existentes para usar o Graviton

Você pode configurar seus aplicativos Amazon EMR Serverless existentes para usar a arquitetura Graviton (`arm64`) com o SDK ou o EMR Studio. AWS CLI

Para converter uma aplicação existente de `x86` para `arm64`

1. Confirme se você está usando a versão principal mais recente da [AWS CLI](#) ou do [SDK](#) com suporte ao parâmetro `architecture`.
2. Confirme se não há trabalhos em execução e interrompa a aplicação.

```
aws emr-serverless stop-application \  
  --application-id application-id \  
  --region us-west-2
```

3. Para atualizar a aplicação e usar o Graviton, especifique `ARM64` para o parâmetro `architecture` na API `update-application`.

```
aws emr-serverless update-application \  
  --application-id application-id \  
  --architecture 'ARM64' \  
  --region us-west-2
```

4. Para verificar se a arquitetura da CPU do aplicativo está agora ARM64, use a `get-application` API.

```
aws emr-serverless get-application \  
--application-id application-id \  
--region us-west-2
```

5. Quando estiver tudo pronto, reinicie a aplicação.

```
aws emr-serverless start-application \  
--application-id application-id \  
--region us-west-2
```

Considerações ao usar o Graviton

Antes de iniciar uma aplicação do EMR Sem Servidor usando `arm64` para suporte ao Graviton, confirme os requisitos a seguir.

Compatibilidade de bibliotecas

Ao selecionar o Graviton (`arm64`) como opção de arquitetura, certifique-se de que pacotes e bibliotecas de terceiros sejam compatíveis com a arquitetura ARM de 64 bits. Para obter informações sobre como empacotar bibliotecas Python em um ambiente virtual Python compatível com a arquitetura selecionada, consulte [Uso de bibliotecas Python com o EMR Sem Servidor](#).

Para saber mais, consulte o repositório [AWS Graviton Getting Started](#) em GitHub. Esse repositório contém recursos essenciais que podem ajudar você a começar a usar o Graviton baseado em ARM.

Simultaneidade e enfileiramento de trabalhos para um aplicação do EMR Sem Servidor

A partir da versão 7.0.0 e posterior do Amazon EMR, especifique o tempo limite da fila de execução de trabalhos e a configuração de simultaneidade da aplicação. Quando você especifica essa configuração, o Amazon EMR Sem Servidor começa colocando seu trabalho em fila e inicia a execução com base na utilização simultânea na aplicação. Por exemplo, se a simultaneidade de execução do seu trabalho for 10, somente dez trabalhos serão executados por vez na aplicação. Os trabalhos restantes ficam na fila até que um dos trabalhos em execução seja encerrado. Se o tempo

limite da fila for atingido mais cedo, seu trabalho expirará. Para obter mais informações, consulte [Estados de execução de trabalhos](#).

Principais benefícios da simultaneidade e do enfileiramento

A simultaneidade e o enfileiramento de trabalhos oferecem os seguintes benefícios quando muitos envios de trabalhos são necessários:

- Isso ajuda a controlar trabalhos de execução simultâneos para usar com eficiência os limites de capacidade no nível da aplicação.
- A fila pode conter uma expansão repentina de envios de trabalhos, com uma definição de tempo limite configurável.

Conceitos básicos de simultaneidade e enfileiramento

Os procedimentos a seguir mostram algumas maneiras diferentes de implementar a simultaneidade e o enfileiramento.

Usando o AWS CLI

1. Crie uma aplicação do Amazon EMR Sem Servidor com tempo limite de fila e execuções de trabalhos simultâneos:

```
aws emr-serverless create-application \  
--release-label emr-7.0.0 \  
--type SPARK \  
--scheduler-configuration '{"maxConcurrentRuns": 1, "queueTimeoutMinutes": 30}'
```

2. Atualize uma aplicação para alterar o tempo limite da fila e a simultaneidade de trabalhos:

```
aws emr-serverless update-application \  
--application-id application-id \  
--scheduler-configuration '{"maxConcurrentRuns": 5, "queueTimeoutMinutes": 30}'
```

Note

Você pode atualizar a aplicação existente para habilitar a simultaneidade e o enfileiramento de trabalhos. Para fazer isso, a aplicação deve ter um rótulo de lançamento emr-7.0.0 ou posterior.

Usando o Console de gerenciamento da AWS

As seguintes etapas mostram como começar a usar a simultaneidade e o enfileiramento de trabalhos usando o Console de gerenciamento da AWS:

1. Acesse EMR Studio e escolha criar uma aplicação com o rótulo de lançamento EMR-7.0.0 ou superior.
2. Em Opções de configuração da aplicação, selecione a opção Usar configurações personalizadas.
3. Em Configurações adicionais, há uma seção para Configurações de execução do trabalho. Selecione a opção Habilitar simultaneidade de trabalhos para habilitar o recurso.
4. Após a seleção, escolha Execuções de trabalhos simultâneos e Tempo limite da fila para configurar o número de execuções de trabalhos simultâneos e o tempo limite da fila, respectivamente. Se você não inserir valores para essas configurações, os valores padrão serão usados.
5. Escolha Criar aplicações e a aplicação será criada com esse recurso habilitado. Para verificar, acesse o painel, selecione sua aplicação e confira na guia de propriedades para determinar se o recurso está habilitado.

Após a configuração, envie trabalhos com esse recurso habilitado.

Considerações sobre simultaneidade e enfileiramento

Leve as seguintes informações em consideração ao implementar a simultaneidade e o enfileiramento:

- A simultaneidade e o enfileiramento de trabalhos têm suporte na versão 7.0.0 e superior do Amazon EMR.

- A simultaneidade e o enfileiramento de tarefas estão habilitados por padrão no Amazon EMR versão 7.3.0 e superior.
- Você não pode atualizar a simultaneidade de uma aplicação no estado STARTED.
- O intervalo válido para `maxConcurrentRuns` é de 1 a 1000 e, para `queueTimeoutMinutes`, ele é de 15 a 720.
- Um máximo de 2.000 trabalhos podem estar no estado QUEUED em uma conta.
- A simultaneidade e o enfileiramento se aplicam a trabalhos em lote e de streaming. Ele não pode ser usado para trabalhos interativos. Para obter mais informações, consulte [Executar workloads interativas com o EMR Sem Servidor através do EMR Studio](#).

Transferência de dados para a classe S3 Express One Zone com o EMR Sem Servidor

Com as versões 7.2.0 e superiores do Amazon EMR, use o EMR Sem Servidor com a classe de armazenamento [Amazon S3 Express One Zone](#) para melhorar a performance ao executar trabalhos e workloads. O S3 Express One Zone é uma classe de armazenamento de zona única e alta performance do Amazon S3 que oferece acesso consistente a dados de milissegundos de um dígito para a maioria das aplicações sensíveis à latência. Na hora da execução, o S3 Express One Zone oferece o armazenamento de objetos na nuvem com a menor latência e a maior performance do Amazon S3.

Pré-requisitos

- Permissões do S3 Express One Zone: quando o S3 Express One Zone inicialmente executa uma ação como GET, LIST ou PUT em um objeto do S3, a classe de armazenamento chama `CreateSession` em seu nome. Sua política do IAM deve conceder a permissão `s3express:CreateSession` para que o conector S3A possa invocar a API `CreateSession`. Para obter um exemplo de política com essa permissão, consulte [Conceitos básicos da classe S3 Express One Zone](#).
- Conector S3A: para configurar o Spark para acessar dados de um bucket do Amazon S3 que usa a classe de armazenamento S3 Express One Zone, use o conector S3A do Apache Hadoop. Para usar o conector, certifique-se de que todos os S3 URIs usem o `s3a` esquema. Caso contrário, altere a implementação do sistema de arquivos usado para os esquemas do `s3` e do `s3n`.

Para alterar o esquema do `s3`, especifique as seguintes configurações de cluster:

```
[
  {
    "Classification": "core-site",
    "Properties": {
      "fs.s3.impl": "org.apache.hadoop.fs.s3a.S3AFileSystem",
      "fs.AbstractFileSystem.s3.impl": "org.apache.hadoop.fs.s3a.S3A"
    }
  }
]
```

Para alterar o esquema do s3n, especifique as seguintes configurações de cluster:

```
[
  {
    "Classification": "core-site",
    "Properties": {
      "fs.s3n.impl": "org.apache.hadoop.fs.s3a.S3AFileSystem",
      "fs.AbstractFileSystem.s3n.impl": "org.apache.hadoop.fs.s3a.S3A"
    }
  }
]
```

Conceitos básicos da classe S3 Express One Zone

Siga estas etapas para começar a usar o S3 Express One Zone.

1. [Crie um endpoint da VPC](#). Adicione o endpoint `com.amazonaws.us-west-2.s3express` ao endpoint da VPC.
2. Siga [Getting started with Amazon EMR Sem Servidor](#) para criar uma aplicação com o rótulo de lançamento 7.2.0 ou superior do Amazon EMR.
3. [Configure a aplicação](#) para usar o endpoint da VPC recém-criado, um grupo de sub-rede privada e um grupo de segurança.
4. Adicione a permissão `CreateSession` ao perfil de execução do trabalho.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Resource": [
        "*"
      ],
      "Action": [
        "s3express:CreateSession"
      ],
      "Sid": "AllowS3EXPRESSCreatesession"
    }
  ]
}
```

```
]
}
```

5. Execute o trabalho. Observe que você deve usar o esquema S3A para acessar os buckets do S3 Express One Zone.

```
aws emr-serverless start-job-run \  
--application-id <application-id> \  
--execution-role-arn <job-role-arn> \  
--name <job-run-name> \  
--job-driver '{  
  "sparkSubmit": {  
  
    "entryPoint": "s3a://<DOC-EXAMPLE-BUCKET>/scripts/wordcount.py",  
    "entryPointArguments":["s3a://<DOC-EXAMPLE-BUCKET>/emr-serverless-spark/output"],  
    "sparkSubmitParameters": "--conf spark.executor.cores=4  
--conf spark.executor.memory=8g --conf spark.driver.cores=4  
--conf spark.driver.memory=8g --conf spark.executor.instances=2  
--conf spark.hadoop.fs.s3a.change.detection.mode=none  
--conf spark.hadoop.fs.s3a.endpoint.region={<AWS_REGION>}  
--conf spark.hadoop.fs.s3a.select.enabled=false  
--conf spark.sql.sources.fastS3PartitionDiscovery.enabled=false  
}'
```

Execução de trabalhos

Após provisionar sua aplicação, envie trabalhos para ela. Esta seção aborda como usar o AWS CLI para executar esses trabalhos. Esta seção também identifica os valores padrão para cada tipo de aplicação disponível no EMR Sem Servidor.

Tópicos

- [Estados de execução de trabalho](#)
- [Cancelamento da execução de trabalhos do EMR Sem Servidor com período de tolerância](#)
- [Execução de trabalhos no console do EMR Studio](#)
- [Executando trabalhos a partir do AWS CLI](#)
- [Política do IAM de execução](#)
- [Uso de discos otimizados para embaralhamento](#)
- [Usando o armazenamento sem servidor para o Amazon EMR Serverless](#)
- [Trabalhos de streaming para processamento contínuo de dados transmitidos](#)
- [Uso das configurações do Spark ao executar trabalhos do EMR Sem Servidor](#)
- [Uso das configurações do Hive ao executar trabalhos do EMR Sem Servidor](#)
- [Resiliência de trabalhos no EMR Sem Servidor](#)
- [Configuração da metastore para EMR Sem Servidor](#)
- [Acessando dados do S3 em outra AWS conta do EMR Serverless](#)
- [Solução de problemas no EMR Sem Servidor](#)
- [Habilitando a alocação de custos no nível de Job](#)

Estados de execução de trabalho

Quando você envia uma execução de trabalho para uma fila de trabalhos do Amazon EMR Sem Servidor, a execução do trabalho entra no estado SUBMITTED. O estado de um trabalho passa de SUBMITTED a RUNNING até atingir FAILED, SUCCESS ou CANCELLING.

As execuções de trabalhos podem ter os seguintes estados:

Estado	Description
Enviado	O estado inicial do trabalho quando você envia uma execução de trabalho ao EMR Sem Servidor. O trabalho espera ser programado para a aplicação. O EMR Sem Servidor começa a priorizar e programar a execução do trabalho.
Enfileirado	A execução do trabalho aguarda nesse estado quando a simultaneidade de execução do trabalho no nível da aplicação está totalmente ocupada. Para obter mais informações sobre enfileiramento e simultaneidade, consulte Simultaneidade e enfileiramento de trabalhos para um aplicação do EMR Sem Servidor .
Pendente	O programador está avaliando a execução do trabalho para priorizar e programar a execução da aplicação.
Programado	O EMR Sem Servidor programou a execução do trabalho para a aplicação e está alocando recursos para executar o trabalho.
Em execução	O EMR Sem Servidor alocou os recursos que o trabalho inicialmente precisa, e o trabalho está sendo executado na aplicação. Em aplicações do Spark, isso significa que o processo do driver do Spark está no estado <code>running</code> .
Com falha	O EMR Sem Servidor falhou ao enviar a execução do trabalho para a aplicação ou o envio foi concluído sem êxito. Consulte <code>StateDetails</code> para obter informações adicionais sobre essa falha de trabalho.

Estado	Description
Bem-sucedida	A execução do trabalho foi concluída com êxito.
Cancelando	A API <code>CancelJobRun</code> solicitou o cancelamento da execução do trabalho ou a execução do trabalho atingiu o tempo limite. O EMR Sem Servidor está tentando cancelar o trabalho na aplicação e liberar os recursos.
Cancelado	A execução do trabalho foi cancelada com êxito e os recursos usados foram liberados.

Cancelamento da execução de trabalhos do EMR Sem Servidor com período de tolerância

Em sistemas de processamento de dados, terminações abruptas podem levar ao desperdício de recursos, operações incompletas e possíveis inconsistências de dados. O Amazon EMR Sem Servidor permite que você especifique um período de tolerância ao cancelar a execução de trabalhos. Esse recurso permite tempo para a limpeza adequada e a conclusão do trabalho em andamento antes do término do trabalho.

Ao cancelar a execução de um trabalho, especifique um período de tolerância (em segundos) usando o parâmetro `shutdownGracePeriodInSeconds` durante o qual o trabalho pode realizar operações de limpeza antes do término final. O comportamento e as configurações padrão variam entre trabalhos em lote e de transmissão.

Período de tolerância para trabalhos em lote

Para trabalhos em lote, o EMR Sem Servidor permite que você implemente operações de limpeza personalizadas que são executadas durante o período de tolerância. Você pode registrar essas operações de limpeza como parte do hook de desligamento da JVM no código da aplicação.

Comportamento padrão

O comportamento padrão para o desligamento é não ter período de tolerância. Consiste nas duas ações a seguir:

- Término imediato
- Os recursos são liberados imediatamente

Opções de configuração

Você pode especificar configurações que resultam em um desligamento normal:

- Intervalo válido para o período de tolerância de desligamento: 15 a 1800 segundos (opcional)
- Término imediato (sem período de tolerância): 0 segundo

Habilitar desligamento normal

Para implementar o desligamento normal para trabalhos em lote, siga estas etapas:

1. Adicione um hook de desligamento no código da sua aplicação contendo uma lógica de desligamento personalizada.

Example in Scala

```
import org.apache.hadoop.util.ShutdownHookManager

// Register shutdown hook with priority (second argument)
// Higher priority hooks run first
ShutdownHookManager.get().addShutdownHook(() => {
  logger.info("Performing cleanup operations...")
}, 100)
```

Usar o [ShutdownHookManager](#)

Example in PySpark

```
import atexit

def cleanup():
    # Your cleanup logic here
    print("Performing cleanup operations...")

# Register the cleanup function
atexit.register(cleanup)
```

2. Especifique um período de tolerância ao cancelar o trabalho para permitir que os hooks adicionados anteriormente sejam executados.

Exemplo

```
# Default (immediate termination)
aws emr-serverless cancel-job-run \
  --application-id APPLICATION_ID \
  --job-run-id JOB_RUN_ID

# With 5-minute grace period
aws emr-serverless cancel-job-run \
  --application-id APPLICATION_ID \
  --job-run-id JOB_RUN_ID \
  --shutdown-grace-period-in-seconds 300
```

Período de tolerância para trabalhos de transmissão

No Spark Structured Streaming, em que os cálculos envolvem leitura ou gravação em fontes de dados externas, desligamentos abruptos podem levar a resultados indesejados. As tarefas de transmissão processam dados em microlotes, e interromper essas operações no meio do caminho pode resultar em processamento duplicado em tentativas subsequentes. Isso acontece quando o ponto de verificação mais recente do microlote anterior não foi gravado, fazendo com que os mesmos dados sejam processados novamente quando o trabalho de transmissão for reiniciado. Esse processamento duplicado não só desperdiça recursos de computação, como também pode afetar as operações comerciais, tornando crucial evitar paralisações abruptas.

O EMR Sem Servidor fornece suporte integrado para um desligamento normal por meio de um receptor de consultas de transmissão. Isso garante a conclusão adequada dos microlotes contínuos antes do término do trabalho. O serviço gerencia automaticamente o desligamento normal entre microlotes para aplicações de transmissão, garantindo que o microlote atual conclua o processamento, que os pontos de verificação sejam gravados adequadamente e que o contexto de transmissão seja encerrado de forma limpa, sem ingerir novos dados durante o processo de desligamento.

Comportamento padrão

- Período de tolerância de 120 segundos habilitado por padrão.
- O receptor de consultas de transmissão integrado gerencia o desligamento normal.

Opções de configuração

- Intervalo válido para o período de tolerância de desligamento: 15 a 1800 segundos (opcional)
- Término imediato: 0 segundo

Habilitar desligamento normal

Para implementar um desligamento normal para trabalhos de transmissão:

Especifique um período de tolerância ao cancelar o trabalho para permitir que o microlote em andamento seja concluído.

Exemplo

```
# Default graceful shutdown (120 seconds)
aws emr-serverless cancel-job-run \
  --application-id APPLICATION_ID \
  --job-run-id JOB_RUN_ID

# Custom grace period (e.g. 300 seconds)
aws emr-serverless cancel-job-run \
  --application-id APPLICATION_ID \
  --job-run-id JOB_RUN_ID \
  --shutdown-grace-period-in-seconds 300

# Immediate Termination
aws emr-serverless cancel-job-run \
  --application-id APPLICATION_ID \
  --job-run-id JOB_RUN_ID \
  --shutdown-grace-period-in-seconds 0
```

Adicione hooks de desligamento personalizados (opcional)

Enquanto o EMR Sem Servidor gerencia o desligamento normal por padrão por meio de seu receptor de consultas de transmissão integrado, você pode, opcionalmente, implementar uma lógica de desligamento personalizada para consultas de transmissão individuais. O EMR Serverless registra seu elegante ouvinte de desligamento com prioridade 60 (usando). ShutdownHookManager Como os hooks de maior prioridade são executados primeiro, você pode registrar suas operações de limpeza personalizadas com uma prioridade maior que 60 para garantir que elas sejam executadas antes do início do processo de desligamento do EMR Sem Servidor.

Para adicionar um hook personalizado, consulte o primeiro exemplo neste tópico que mostra como adicionar um hook de desligamento no código da aplicação. Aqui, 100 é a prioridade, que é maior que 60. Portanto, esse hook de desligamento funciona primeiro.

Note

Hooks de desligamento personalizados são opcionais e não são necessários para uma funcionalidade de desligamento normal, que é gerenciada automaticamente pelo EMR Sem Servidor.

Cobranças do período de tolerância e duração do lote

Se o valor padrão para o período de tolerância (120 segundos) for usado:

- Se a duração do lote for inferior a 120 segundos, apenas haverá cobranças pelo tempo real necessário para concluir o lote.
- Se a duração do lote exceder 120 segundos, a cobrança ocorrerá com base no período máximo de tolerância (120 segundos), mas a consulta poderá não ser encerrada normalmente, pois será encerrada à força.

Para otimizar os custos e garantir um desligamento adequado:

- Para durações de lote superiores a 120 segundos: considere aumentar o período de tolerância para corresponder à duração do seu lote
- Para lotes com duração inferior a 120 segundos: não é necessário ajustar o período de tolerância, pois você só será cobrado pelo tempo real de processamento

Considerações

Comportamento durante o período de tolerância

- O período de tolerância fornece tempo para que seus hooks de desligamento registrados sejam concluídos.
- O trabalho termina assim que o processo de desligamento termina, mesmo que seja bem antes do período de tolerância.
- Se as operações de limpeza excederem o período de tolerância, o trabalho será encerrado à força.

Comportamento do serviço

- O encerramento do período de tolerância apenas está disponível para trabalhos no estado RUNNING.
- Solicitações de cancelamento subsequentes durante o estado CANCELLING serão ignoradas.
- Se o EMR Sem Servidor falhar em iniciar o desligamento do período de tolerância devido a erros internos de serviço:
 - O serviço será repetido por até 2 minutos.
 - Se as novas tentativas não tiverem êxito, o trabalho será encerrado à força.

Faturamento

Os trabalhos são cobrados pelos recursos computacionais usados até o encerramento completo do trabalho, incluindo qualquer tempo gasto durante o período de tolerância.

Execução de trabalhos no console do EMR Studio

Você pode enviar execuções de trabalhos para aplicações do EMR Sem Servidor e acessar os trabalhos no console do EMR Studio. Para criar ou navegar até a aplicação do EMR Sem Servidor no console do EMR Studio, siga as instruções em [Getting started from the console](#).

Enviar uma tarefa

Na página Enviar trabalho, envie um trabalho a uma aplicação do EMR Sem Servidor da forma mostrada a seguir.

Spark

1. No campo Nome, insira um nome para a execução do trabalho.
2. No campo Perfil de runtime, insira o nome do perfil do IAM que a aplicação do EMR Sem Servidor pode assumir para a execução do trabalho. Para saber mais sobre perfis de runtime, consulte [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#).
3. No campo Localização do script, insira a localização do Amazon S3 do script ou JAR que você deseja executar. Em trabalhos do Spark, o script pode ser um arquivo Python (.py) ou um arquivo JAR (.jar).
4. Se a localização do script for um arquivo JAR, insira o nome da classe que é o ponto de entrada do trabalho no campo Classe principal.

5. (Opcional) Insira valores nos campos restantes.

- Argumentos do script: insira os argumentos que você deseja passar ao script JAR ou Python principal. Seu código lê esses parâmetros. Separe cada argumento da matriz por uma vírgula.
- Propriedades do Spark: expanda a seção de propriedades do Spark e insira quaisquer parâmetros de configuração do Spark nesse campo.

Note

Se você especificar os tamanhos do driver e do executor do Spark, considere a sobrecarga de memória. Especifique os valores de sobrecarga de memória nas propriedades `spark.driver.memoryOverhead` e `spark.executor.memoryOverhead`. A sobrecarga de memória tem um valor padrão de 10% da memória do contêiner, com um mínimo de 384 MB. Juntas, a memória do executor e a sobrecarga de memória não podem exceder a memória do trabalhador. Por exemplo, a `spark.executor.memory` máxima em um trabalhador de 30 GB deve ser 27 GB.

- Configuração de trabalho: especifique qualquer configuração de trabalho nesse campo. Você pode usar essas configurações de trabalho para substituir as configurações padrão de aplicações.
- Configurações adicionais: ative ou desative o AWS Glue Data Catalog como uma metastore e modifique as configurações de log da aplicação. Para saber mais sobre configurações do metastore, consulte [Configuração da metastore para EMR Sem Servidor](#). Para saber mais sobre as opções de registro em log de aplicações, consulte [Armazenamento de logs](#).
- Tags: atribua tags personalizadas à aplicação.

6. Escolha Enviar trabalho.

Hive

1. No campo Nome, insira um nome para a execução do trabalho.
2. No campo Perfil de runtime, insira o nome do perfil do IAM que a aplicação do EMR Sem Servidor pode assumir para a execução do trabalho.

3. No campo Localização do script, insira a localização do Amazon S3 do script ou JAR que você deseja executar. Em trabalhos do Hive, o script deve ser um arquivo do Hive (.sql).
4. (Opcional) Insira valores nos campos restantes.
 - Localização do script de inicialização: insira a localização do script que inicializa as tabelas antes da execução do script do Hive.
 - Propriedades do Hive: expanda a seção Propriedades do Hive e insira quaisquer parâmetros de configuração do Hive nesse campo.
 - Configuração de trabalho: especifique qualquer configuração de trabalho. Você pode usar essas configurações de trabalho para substituir as configurações padrão de aplicações. Em trabalhos do Hive, `hive.exec.scratchdir` e `hive.metastore.warehouse.dir` são propriedades obrigatórias na configuração `hive-site`.

```
{
  "applicationConfiguration": [
    {
      "classification": "hive-site",
      "configurations": [],
      "properties": {
        "hive.exec.scratchdir": "s3://DOC-EXAMPLE_BUCKET/hive/scratch",
        "hive.metastore.warehouse.dir": "s3://DOC-EXAMPLE_BUCKET/hive/warehouse"
      }
    }
  ],
  "monitoringConfiguration": {}
}
```

- Configurações adicionais — Ative ou desative o AWS Glue Data Catalog como um metastore e modifique as configurações de registro do aplicativo. Para saber mais sobre configurações do metastore, consulte [Configuração da metastore para EMR Sem Servidor](#). Para saber mais sobre as opções de registro em log de aplicações, consulte [Armazenamento de logs](#).
 - Tags: atribua quaisquer tags personalizadas à aplicação.
5. Escolha Enviar trabalho.

Acessar execuções de trabalhos

Na guia Execuções de trabalhos na página Detalhes de uma aplicação, acesse execuções de trabalhos e realize as ações a seguir para execuções de trabalhos.

Cancelar trabalho: para cancelar a execução de um trabalho que esteja no estado RUNNING, escolha essa opção. Para saber mais sobre as transições de execução de trabalhos, consulte [Estados de execução de trabalho](#).

Clonar trabalho: para clonar uma execução de trabalho anterior e reenviá-la, escolha essa opção.

Executando trabalhos a partir do AWS CLI

Você pode criar, descrever e excluir trabalhos individuais na AWS CLI. Você também pode listar todos os seus trabalhos para acessá-los rapidamente.

Para enviar um novo trabalho, use `start-job-run`. Forneça o ID da aplicação que deseja executar, com as propriedades específicas do trabalho. Para exemplos do Spark, consulte [Uso das configurações do Spark ao executar trabalhos do EMR Sem Servidor](#). Para exemplos do Hive, consulte [Uso das configurações do Hive ao executar trabalhos do EMR Sem Servidor](#). Esse comando retorna o `application-id`, o ARN e o novo `job-id`.

Cada execução de trabalho tem uma duração de tempo limite definida. Se a execução do trabalho exceder essa duração, o EMR Sem Servidor o cancelará automaticamente. O tempo limite padrão é de 12 horas. Ao iniciar a execução do trabalho, defina essa configuração de tempo limite para um valor que atenda aos requisitos do trabalho. Configure o valor com a propriedade `executionTimeoutMinutes`.

```
aws emr-serverless start-job-run \  
  --application-id application-id \  
  --execution-role-arn job-role-arn \  
  --execution-timeout-minutes 15 \  
  --job-driver '{  
    "hive": {  
      "query": "s3://amzn-s3-demo-bucket/scripts/create_table.sql",  
      "parameters": "--hiveconf hive.exec.scratchdir=s3://amzn-s3-demo-bucket/  
hive/scratch --hiveconf hive.metastore.warehouse.dir=s3://amzn-s3-demo-bucket/hive/  
warehouse"  
    }  
  }' \  

```

```
--configuration-overrides '{
  "applicationConfiguration": [{
    "classification": "hive-site",
    "properties": {
      "hive.client.cores": "2",
      "hive.client.memory": "4GIB"
    }
  }]
}'
```

Para descrever um trabalho, use `get-job-run`. Esse comando retorna as configurações específicas do trabalho e a capacidade definida do novo trabalho.

```
aws emr-serverless get-job-run \
--job-run-id job-id \
--application-id application-id
```

Para listar seus trabalhos, use `list-job-runs`. Esse comando retorna um conjunto abreviado de propriedades que inclui tipo de trabalho, estado e outros atributos de alto nível. Se você não quiser acessar todos os trabalhos, especifique o número máximo de trabalhos que deseja acessar, até 50. O exemplo a seguir especifica que você deseja acessar seus dois últimos trabalhos executados.

```
aws emr-serverless list-job-runs \
--max-results 2 \
--application-id application-id
```

Use `cancel-job-run` para cancelar um trabalho. Forneça o `application-id` e o `job-id` do trabalho que você deseja cancelar.

```
aws emr-serverless cancel-job-run \
--job-run-id job-id \
--application-id application-id
```

Para obter mais informações sobre como executar trabalhos a partir do AWS CLI, consulte a Referência da API do [EMR Serverless](#).

Política do IAM de execução

É possível especificar uma política de execução do IAM, além de um perfil de execução, ao enviar execuções de trabalhos no EMR Sem Servidor. As permissões resultantes assumidas pelo

trabalho são a interseção das permissões no perfil de execução e na Política do IAM de execução especificada.

Conceitos básicos

Etapas para usar a política do IAM de execução:

Crie um aplicação `emr-serverless` ou use uma existente e, em seguida, execute o seguinte comando da AWS CLI para iniciar a execução de um trabalho com uma política do IAM inline:

```
aws emr-serverless start-job-run --region us-west-2 \  
  --application-id application-id \  
  --execution-role-arn execution-role-arn \  
  --job-driver job-driver-options \  
  --execution-iam-policy '{"policy": "inline-policy"}'
```

Exemplos de comandos da CLI

Se tivermos a seguinte política armazenada no arquivo `policy.json` na máquina:

JSON

```
{  
  "Version": "2012-10-17",  
  "Statement": [  
    {  
      "Effect": "Allow",  
      "Action": [  
        "s3:GetObject",  
        "s3:ListBucket"  
      ],  
      "Resource": [  
        "arn:aws:s3::my-test-bucket",  
        "arn:aws:s3::my-test-bucket/*"  
      ],  
      "Sid": "AllowS3GetObject"  
    }  
  ]  
}
```

Em seguida, podemos iniciar um trabalho com essa política usando o seguinte AWS CLI comando:

```
aws emr-serverless start-job-run --region us-west-2 \  
  --application-id application-id \  
  --execution-role-arn execution-role-arn \  
  --job-driver job-driver-options \  
  --execution-iam-policy '{  
    "policy": '$(jq -c '. | @json' policy.json)'  
  }'
```

Você também pode usar ambas as políticas AWS e as gerenciadas pelo cliente, especificando-as ARNs por meio de:

```
aws emr-serverless start-job-run --region us-west-2 \  
  --application-id application-id \  
  --execution-role-arn execution-role-arn \  
  --job-driver job-driver-options \  
  --execution-iam-policy '{  
    "policyArns": [  
      "arn:aws:iam::aws:policy/AmazonS3FullAccess",  
      "arn:aws:iam::aws:policy/CloudWatchLogsFullAccess"  
    ]  
  }'
```

Também é possível especificar uma política de IAM em linha e uma política gerenciada ARNs na mesma solicitação:

```
aws emr-serverless start-job-run --region us-west-2 \  
  --application-id application-id \  
  --execution-role-arn execution-role-arn \  
  --job-driver job-driver-options \  
  --execution-iam-policy '{  
    "policy": '$(jq -c '. | @json' policy.json)',  
    "policyArns": [  
      "arn:aws:iam::aws:policy/AmazonS3FullAccess",  
      "arn:aws:iam::aws:policy/CloudWatchLogsFullAccess"  
    ]  
  }'
```

Observações importantes

- O campo `policy` de `execution-role-policy` pode ter um comprimento máximo de 2048 caracteres.
- A string de política do IAM inline especificada no campo `policy` do `execution-iam-policy` deve estar em conformidade com o padrão de string json, sem novas linhas e com aspas escapadas, como no exemplo anterior.
- Uma lista de até 10 políticas gerenciadas ARNs pode ser especificada como um valor para `execution-iam-policy` o `policyArns` campo `s`.
- A política gerenciada ARNs deve ser uma lista de ARN válidos AWS ou gerenciados pelo cliente. Quando um ARN de política gerenciada pelo cliente é especificado, a política deve pertencer à mesma AWS conta do EMR-S. JobRun
- Quando a política do IAM inline e políticas gerenciadas são usadas, o texto simples que você usa para as políticas inline e gerenciadas combinadas não pode exceder 2.048 caracteres.
- As permissões resultantes assumidas pelo JobRun são a interseção das permissões na função de execução e na política de execução IAM especificada.

Interseção de políticas

As permissões resultantes assumidas pelo trabalho são a interseção das permissões no perfil de execução e na Política do IAM de execução especificada. O que significa que qualquer permissão necessária deverá ser especificada em ambos os lugares JobRun para funcionar. Porém, é possível especificar uma declaração de permissão geral adicional na política inline para quaisquer permissões que você não pretende atualizar ou substituir.

Exemplo

Dada a seguinte política de perfil do IAM de execução:

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
```

```

    "Action": [
      "s3:*"
    ],
    "Resource": [
      "*"
    ],
    "Sid": "AllowS3"
  },
  {
    "Effect": "Allow",
    "Action": [
      "logs:DescribeLogGroups"
    ],
    "Resource": [
      "arn:aws:logs:us-west-2:123456789012:log-group::log-stream"
    ],
    "Sid": "AllowLOGSDescribeLogGroups"
  },
  {
    "Effect": "Allow",
    "Action": [
      "dynamodb:DescribeTable"
    ],
    "Resource": [
      "arn:aws:dynamodb:*:*:table/MyCompany1table"
    ],
    "Sid": "AllowDYNAMODescribeTable"
  }
]
}

```

E a seguinte política do IAM embutida:

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",

```

```

    "s3:ListBucket"
  ],
  "Resource": [
    "arn:aws:s3:::my-test-bucket/tenant1",
    "arn:aws:s3:::my-test-bucket/tenant1/*"
  ],
  "Sid": "AllowS3GetObject"
},
{
  "Effect": "Allow",
  "Action": [
    "logs:*",
    "dynamodb:*"
  ],
  "Resource": [
    "*"
  ],
  "Sid": "AllowLOGS"
}
]
}

```

As permissões resultantes assumidas pelo JobRun são:

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3:::my-test-bucket/tenant1",
        "arn:aws:s3:::my-test-bucket/tenant1/*"
      ],
      "Sid": "AllowS3GetObject"
    },
    {

```

```
"Effect": "Allow",
"Action": [
  "logs:DescribeLogGroups"
],
"Resource": [
  "arn:aws:logs:us-west-2:123456789012:log-group:log-stream"
],
"Sid": "AllowLOGSDescribeLogGroups"
},
{
  "Effect": "Allow",
  "Action": [
    "dynamodb:DescribeTable"
  ],
  "Resource": [
    "arn:aws:dynamodb:*:*:table/MyCompany1table"
  ],
  "Sid": "AllowDYNAMODBDescribeTable"
}
]
```

Uso de discos otimizados para embaralhamento

Com as versões 7.1.0 e superiores do Amazon EMR, use discos otimizados para reprodução aleatória ao executar trabalhos do Apache Spark ou do Hive para melhorar o desempenho das operações (por segundo) para acelerar a movimentação de dados e reduzir a latência durante I/O-intensive workloads. Compared to standard disks, shuffle-optimized disks provide higher IOPS (I/O as operações aleatórias. Como os discos otimizados para embaralhamento permitem anexar tamanhos de disco de até 2 TB por trabalhador, configure a capacidade apropriada para seus requisitos de workload.

Benefícios principais

Os discos otimizados para embaralhamento fornecem os benefícios a seguir.

- Alto desempenho de IOPS: discos otimizados para embaralhamento fornecem IOPS mais altas do que os discos padrão, levando a um embaralhamento de dados mais eficiente e rápido durante trabalhos do Spark e do Hive e outras workloads com uso intenso de embaralhamento.

- Tamanho de disco maior: como os discos otimizados para embaralhamento oferecem suporte a tamanhos de disco de 20 GB a 2 TB por trabalhador, escolha a capacidade adequada com base nas suas workloads.

Introdução

Confira as etapas a seguir para usar discos otimizados para embaralhamento nos fluxos de trabalho.

Spark

1. Crie uma aplicação do EMR Sem Servidor versão 7.1.0 com o comando a seguir.

```
aws emr-serverless create-application \  
  --type "SPARK" \  
  --name my-application-name \  
  --release-label emr-7.1.0 \  
  --region <AWS_REGION>
```

2. Configure sua tarefa do Spark para incluir os parâmetros a serem executados com `spark.emr-serverless.driver.disk.type` and/or `spark.emr-serverless.executor.disk.type` discos otimizados para reprodução aleatória. Você pode usar um ou ambos os parâmetros, dependendo do seu caso de uso.

```
aws emr-serverless start-job-run \  
  --application-id application-id \  
  --execution-role-arn job-role-arn \  
  --job-driver '{  
    "sparkSubmit": {  
      "entryPoint": "/usr/lib/spark/examples/jars/spark-examples.jar",  
      "entryPointArguments": ["1"],  
      "sparkSubmitParameters": "--class org.apache.spark.examples.SparkPi  
      --conf spark.executor.cores=4  
      --conf spark.executor.memory=20g  
      --conf spark.driver.cores=4  
      --conf spark.driver.memory=8g  
      --conf spark.executor.instances=1  
      --conf spark.emr-serverless.executor.disk.type=shuffle_optimized"  
    }  
  }'
```

Para obter mais informações, consulte [Propriedades de trabalho do Spark](#).

Hive

1. Crie uma aplicação do EMR Sem Servidor versão 7.1.0 com o comando a seguir.

```
aws emr-serverless create-application \  
  --type "HIVE" \  
  --name my-application-name \  
  --release-label emr-7.1.0 \  
  --region <AWS_REGION>
```

2. Configure sua tarefa do Hive para incluir os parâmetros a serem executados com `hive.driver.disk.type` and/or `hive.tez.disk.type` discos otimizados para reprodução aleatória. Você pode usar um ou ambos os parâmetros, dependendo do seu caso de uso.

```
aws emr-serverless start-job-run \  
  --application-id application-id \  
  --execution-role-arn job-role-arn \  
  --job-driver '{  
    "hive": {  
      "query": "s3://<DOC-EXAMPLE-BUCKET>/emr-serverless-hive/query/hive-  
query.q1",  
      "parameters": "--hiveconf hive.log.explain.output=false"  
    }  
  }' \  
  --configuration-overrides '{  
    "applicationConfiguration": [{  
      "classification": "hive-site",  
      "properties": {  
        "hive.exec.scratchdir": "s3://<DOC-EXAMPLE-BUCKET>/emr-  
serverless-hive/hive/scratch",  
        "hive.metastore.warehouse.dir": "s3://<DOC-EXAMPLE-BUCKET>/emr-  
serverless-hive/hive/warehouse",  
        "hive.driver.cores": "2",  
        "hive.driver.memory": "4g",  
        "hive.tez.container.size": "4096",  
        "hive.tez.cpu.vcores": "1",  
        "hive.driver.disk.type": "shuffle_optimized",  
        "hive.tez.disk.type": "shuffle_optimized"  
      }  
    }  
  ]'  
}
```

Para obter mais informações, confira [Hive job properties](#).

Configuração de aplicações com capacidade pré-inicializada

Confira os exemplos a seguir para criar aplicações com base na versão 7.1.0 do Amazon EMR. Essas aplicações têm as seguintes propriedades:

- 5 drivers do Spark pré-inicializados, cada um com 2 vCPUs, 4 GB de memória e 50 GB de disco otimizado para embaralhamento.
- 50 executores pré-inicializados, cada um com 4 vCPUs, 8 GB de memória e 500 GB de disco otimizado para embaralhamento.

Quando essa aplicação executa trabalhos do Spark, ela primeiro consome os trabalhadores pré-inicializados e depois escala os trabalhadores sob demanda até a capacidade máxima de 400 vCPUs e 1.024 GB de memória. Opcionalmente, você pode omitir a capacidade do DRIVER ou EXECUTOR.

Spark

```
aws emr-serverless create-application \  
  --type "SPARK" \  
  --name <my-application-name> \  
  --release-label emr-7.1.0 \  
  --initial-capacity '{  
    "DRIVER": {  
      "workerCount": 5,  
      "workerConfiguration": {  
        "cpu": "2vCPU",  
        "memory": "4GB",  
        "disk": "50GB",  
        "diskType": "SHUFFLE_OPTIMIZED"  
      }  
    },  
    "EXECUTOR": {  
      "workerCount": 50,  
      "workerConfiguration": {  
        "cpu": "4vCPU",  
        "memory": "8GB",  
        "disk": "500GB",  
        "diskType": "SHUFFLE_OPTIMIZED"  
      }  
    }  
  }'
```

```
    }  
  }  
}' \  
--maximum-capacity '{  
  "cpu": "400vCPU",  
  "memory": "1024GB"  
}'
```

Hive

```
aws emr-serverless create-application \  
--type "HIVE" \  
--name <my-application-name> \  
--release-label emr-7.1.0 \  
--initial-capacity '{  
  "DRIVER": {  
    "workerCount": 5,  
    "workerConfiguration": {  
      "cpu": "2vCPU",  
      "memory": "4GB",  
      "disk": "50GB",  
      "diskType": "SHUFFLE_OPTIMIZED"  
    }  
  },  
  "EXECUTOR": {  
    "workerCount": 50,  
    "workerConfiguration": {  
      "cpu": "4vCPU",  
      "memory": "8GB",  
      "disk": "500GB",  
      "diskType": "SHUFFLE_OPTIMIZED"  
    }  
  }  
}' \  
--maximum-capacity '{  
  "cpu": "400vCPU",  
  "memory": "1024GB"  
}'
```

Usando o armazenamento sem servidor para o Amazon EMR Serverless

Com as versões 7.12 e superiores do Amazon EMR, use armazenamento sem servidor ao executar trabalhos do Apache Spark para eliminar o provisionamento de disco local e reduzir os custos de processamento de dados, além de evitar falhas de trabalho causadas por restrições de capacidade de disco. O armazenamento sem servidor gerencia automaticamente as operações de reprodução aleatória, vazamento de disco e armazenamento em cache de disco para suas tarefas sem exigir configuração de capacidade e armazena dados intermediários sem nenhum custo. O Amazon EMR Serverless armazena dados intermediários em um armazenamento sem servidor totalmente gerenciado que é escalado automaticamente com base nas demandas da carga de trabalho e permite que o Spark libere funcionários da computação imediatamente quando ociosos, reduzindo os custos de computação.

Benefícios principais

O armazenamento sem servidor para EMR Serverless oferece os seguintes benefícios.

- **Armazenamento com configuração zero** — O armazenamento sem servidor elimina a necessidade de configurar o tipo e o tamanho do disco local para cada aplicativo ou tarefa. O EMR Serverless gerencia automaticamente as operações de dados intermediárias sem planejamento de capacidade.
- **Evita falhas de trabalho por meio do escalonamento automático** — A capacidade de armazenamento é dimensionada automaticamente com base na demanda da carga de trabalho, evitando falhas de trabalho devido à capacidade insuficiente do disco.
- **Custos reduzidos de processamento de dados** — O armazenamento sem servidor reduz os custos de processamento por meio de dois mecanismos. Primeiro, o armazenamento intermediário de dados é fornecido sem nenhum custo — você paga somente pelos recursos de computação e memória. Em segundo lugar, o armazenamento desacoplado com a alocação dinâmica de recursos do Spark permite que o Spark libere trabalhadores imediatamente quando estão ociosos, em vez de retê-los para preservar dados intermediários em discos locais. Isso permite uma expansão horizontal e ampliação mais rápidas por estágio do Spark, reduzindo os custos de computação para trabalhos em que os estágios posteriores precisam de menos trabalhadores do que os estágios iniciais.

- Armazenamento criptografado com isolamento em nível de trabalho — Todos os dados intermediários são criptografados em trânsito e em repouso com isolamento estrito em nível de trabalho.
- Suporte de controle de acesso refinado — O armazenamento sem servidor oferece suporte ao controle de acesso refinado por meio da integração com o AWS Lake Formation.

Introdução

Veja as etapas a seguir para usar o armazenamento sem servidor para o EMR Serverless em seus fluxos de trabalho do Spark.

1. Crie um aplicativo EMR Serverless

Crie um aplicativo EMR Serverless versão 7.12 (ou posterior) com armazenamento sem servidor ativado definindo a propriedade spark como true na classificação spark-defaults.

spark.aws.serverlessStorage.enabled

```
aws emr-serverless create-application \  
  --type "SPARK" \  
  --name my-application \  
  --release-label emr-7.12.0 \  
  --runtime-configuration '[\  
    "classification": "spark-defaults",\  
    "properties": {\  
      "spark.aws.serverlessStorage.enabled": "true"  
    }\  
  ]]' \  
  --region <AWS_REGION>
```

2. Comece um trabalho no Spark

Inicie a execução de um trabalho em seu aplicativo. Armazenamento sem servidor para EMR O Serverless manipula automaticamente a operação intermediária de dados, como o shuffle, para seu trabalho.

```
aws emr-serverless start-job-run \  
  --application-id <application-id> \  
  --execution-role-arn <job-role-arn> \  
  --job-driver '{  
    "sparkSubmit": {
```

```

    "entryPoint": "s3://<bucket>/script.py",
    "sparkSubmitParameters": "--conf spark.executor.cores=4
      --conf spark.executor.memory=20g
      --conf spark.driver.cores=4
      --conf spark.driver.memory=8g
      --conf spark.executor.instances=10"
  }
}'

```

Você também pode habilitar o armazenamento sem servidor para o EMR Serverless no nível do trabalho, mesmo quando não estiver habilitado no nível do aplicativo. Isso iniciará nós de trabalho habilitados com armazenamento sem servidor para processar seus trabalhos. Você também pode desativar o armazenamento sem servidor para um trabalho específico definindo a mesma propriedade **spark.aws.serverlessStorage.enabled** do Spark como `false`.

```

# Turn on serverless storage for EMR serverless for a specific job
aws emr-serverless start-job-run \
  --application-id <application-id> \
  --execution-role-arn <job-role-arn> \
  --job-driver '{
"sparkSubmit": {
"entryPoint": "/usr/lib/spark/examples/jars/spark-examples.jar",
  "entryPointArguments": ["1"],
  "sparkSubmitParameters": "--class org.apache.spark.examples.SparkPi
    --conf spark.aws.serverlessStorage.enabled": "true"
  }
}'

```

Note

Para continuar usando o provisionamento de disco local tradicional, omita a **spark.aws.serverlessStorage.enabled** configuração ou defina-a como `false`.

Considerações e limitações

- Versão de lançamento — O armazenamento sem servidor é compatível com a versão 7.12 e posterior do Amazon EMR.

- Limites de volume de dados — Cada trabalho pode ler e gravar até um total de 200 GB de dados intermediários por execução do trabalho. Os trabalhos que excederem esse limite falharão com uma mensagem de erro indicando que o limite de armazenamento sem servidor foi atingido.
- Tempo limite de execução do trabalho — O armazenamento sem servidor oferece suporte a trabalhos com tempos limite de execução de até 24 horas. Os trabalhos configurados para tempos limite de execução mais longos falharão com uma mensagem de erro.
- Capacidade pré-inicializada — os trabalhadores com capacidade pré-inicializada não oferecem suporte ao armazenamento sem servidor. Quando você configura a capacidade pré-inicializada, ela só será utilizada por trabalhos que desabilitem explicitamente o armazenamento sem servidor no nível do trabalho. As tarefas com armazenamento sem servidor habilitado sempre provisionarão novos trabalhadores sob demanda e não usarão nenhuma capacidade pré-inicializada, independentemente da configuração no nível do aplicativo.
- Tipos de carga de trabalho — O armazenamento sem servidor não é compatível com trabalhos interativos e de streaming.
- Configuração do trabalhador — O armazenamento sem servidor não é suportado para trabalhadores com 1 ou 2 v. CPUs

Suportado Regiões da AWS

O EMR Serverless oferece suporte ao armazenamento sem servidor nas seguintes regiões:

- Leste dos EUA (N. da Virgínia)
- Oeste dos EUA (Oregon)
- Europa (Irlanda)

Trabalhos de streaming para processamento contínuo de dados transmitidos

Um trabalho de streaming no EMR Sem Servidor é um modo de trabalho que permite analisar e processar dados de streaming quase em tempo real. Esses trabalhos de longa duração sondam dados de streaming e processam continuamente os resultados à medida que os dados chegam. Os trabalhos de streaming são mais adequados para tarefas que exigem processamento de dados em tempo real, como analytics quase em tempo real, detecção de fraudes e mecanismos de recomendações. Os trabalhos de streaming do EMR Sem Servidor fornecem otimizações, como

resiliência de trabalho integrada, monitoramento em tempo real, gerenciamento aprimorado de logs e integração com conectores de streaming.

Estes são alguns casos de uso com trabalhos de streaming:

- **Analytics quase em tempo real:** trabalhos de streaming no Amazon EMR Sem Servidor permitem processar dados de streaming quase em tempo real, para que você possa realizar analytics em tempo real em fluxos de dados contínuos, como dados de log, dados de sensores ou dados de clickstream, a fim de obter insights e tomar decisões oportunas com base nas informações mais recentes.
- **Detecção de fraudes:** use trabalhos de streaming para executar a detecção de fraudes quase em tempo real em transações financeiras, operações de cartão de crédito ou atividades on-line ao analisar fluxos de dados e identificar padrões ou anomalias suspeitas à medida que ocorrem.
- **Mecanismos de recomendação:** os trabalhos de streaming podem processar dados de atividade do usuário e atualizar modelos de recomendações. Isso abre possibilidades de recomendações personalizadas e em tempo real com base em comportamentos e preferências.
- **Analytics de mídia social:** os trabalhos de streaming podem processar dados de mídia social, como tweets, comentários e publicações, possibilitando que as organizações monitorem tendências, analisem sentimentos e gerenciem a reputação da marca quase em tempo real.
- **Analytics da Internet das Coisas (IoT):** trabalhos de streaming podem analisar e lidar com fluxos de dados de alta velocidade de dispositivos de IoT, sensores e máquinas conectadas, para que você possa realizar detecção de anomalias, manutenção preditiva e outros casos de uso de analytics de IoT.
- **Análise de clickstream:** os trabalhos de streaming podem processar e analisar dados de clickstream de sites ou de aplicações móveis. As empresas que usam esses dados podem realizar analytics e saber mais sobre o comportamento do usuário, personalizar as experiências do usuário e otimizar campanhas de marketing.
- **Monitoramento e análise de log:** os trabalhos de streaming também podem processar dados de log de servidores, aplicações e dispositivos de rede. Isso possibilita detecção de anomalias, solução de problemas e integridade e performance do sistema.

Benefícios principais

Os trabalhos de streaming no EMR Sem Servidor fornecem automaticamente resiliência aos trabalhos, que é uma combinação dos seguintes fatores:

- Nova tentativa automática: o EMR Sem Servidor repete automaticamente todos os trabalhos que falharam sem nenhuma entrada manual de sua parte.
- Resiliência da zona de disponibilidade (AZ): o EMR Sem Servidor muda automaticamente os trabalhos de streaming para uma AZ íntegra se a AZ original apresentar problemas.
- Gerenciamento de logs:
 - Alternância de logs: para um gerenciamento mais eficiente do armazenamento em disco, o EMR Sem Servidor alterna periodicamente os logs para trabalhos de streaming longos. Isso evita o acúmulo de logs que podem consumir todo o espaço em disco.
 - Compactação de logs: ajuda a gerenciar e otimizar com eficiência os arquivos de log em persistência gerenciada. A compactação também melhora a experiência de depuração quando você usa o servidor gerenciado de histórico do Spark.

Fontes de dados e coletores de dados compatíveis

O EMR Sem Servidor funciona com várias fontes de dados de entrada e coletores de dados de saída:

- Fontes de dados de entrada compatíveis: Amazon Kinesis Data Streams, Amazon Managed Streaming para Apache Kafka e clusters autogerenciados do Apache Kafka. Por padrão, as versões 7.1.0 e superiores do Amazon EMR incluem o [conector do Amazon Kinesis Data Streams](#), então você não precisa criar ou baixar nenhum pacote adicional.
- Coletores de dados de saída compatíveis — tabelas do AWS Glue Data Catalog, Amazon S3, Amazon Redshift, MySQL, PostgreSQL Oracle, Oracle, Microsoft SQL, Apache Iceberg, Delta Lake e Apache Hudi.

Considerações e limitações

Ao usar trabalhos de streaming, lembre-se das considerações e limitações a seguir.

- Os trabalhos de streaming são compatíveis com as [versões 7.1.0 e superiores do Amazon EMR](#).
- O EMR Sem Servidor espera que os trabalhos de streaming sejam executados por muito tempo, então você não pode definir o tempo limite de execução para limitar o runtime do trabalho.
- Os trabalhos de streaming são compatíveis apenas com o mecanismo do Spark, que é criado sobre o [framework de streaming estruturado](#).

- O EMR Sem Servidor repete indefinidamente os trabalhos de streaming, e você não pode personalizar o número máximo de tentativas. A prevenção contra thrash é incluída automaticamente para interromper a repetição do trabalho se a quantidade de tentativas malsucedidas ultrapassar o limite definido em uma janela por hora. O limite padrão é de cinco tentativas malsucedidas em uma hora. Você pode configurar esse limite entre 1 e 10 tentativas. Para obter mais informações, consulte [Resiliência de trabalhos](#).
- Os trabalhos de streaming têm pontos de verificação para salvar o estado e o progresso do runtime, para que o EMR Sem Servidor possa retomar o trabalho de streaming a partir do ponto de verificação mais recente. Para obter mais informações, consulte [Recuperação de falhas com Checkpointing](#) na documentação do Apache Spark.

Conceitos básicos de trabalhos de streaming

Confira as instruções a seguir para saber como começar a usar trabalhos de streaming.

1. Leia [Getting started with Amazon EMR Sem Servidor to create an application](#). Observe que a aplicação deve executar o [Amazon EMR versão 7.1.0](#) ou superior.
2. Quando seu aplicativo estiver pronto, defina o mode parâmetro STREAMING para enviar um trabalho de streaming, semelhante ao AWS CLI exemplo a seguir.

```
aws emr-serverless start-job-run \  
--application-id <APPLICATION_ID> \  
--execution-role-arn <JOB_EXECUTION_ROLE> \  
--mode 'STREAMING' \  
--job-driver '{  
  "sparkSubmit": {  
    "entryPoint": "s3://<streaming script>",  
    "entryPointArguments": ["s3://<DOC-EXAMPLE-BUCKET-OUTPUT>/output"],  
    "sparkSubmitParameters": "--conf spark.executor.cores=4  
      --conf spark.executor.memory=16g  
      --conf spark.driver.cores=4  
      --conf spark.driver.memory=16g  
      --conf spark.executor.instances=3"  
  }  
'
```

Conectores de streaming compatíveis

Os conectores de streaming facilitam a leitura de dados provenientes de uma fonte de streaming e podem gravar dados em um coletor de streaming.

Estes são os conectores de streaming compatíveis:

Conector do Amazon Kinesis Data Streams

O [conector do Amazon Kinesis Data Streams](#) para Apache Spark permite criar aplicações e pipelines de streaming que consomem dados e gravam dados no Amazon Kinesis Data Streams. O conector comporta um consumo aprimorado de fan-out com um throughput de leitura dedicado de até 2 MB/segundo por fragmento. Por padrão, o Amazon EMR Sem Servidor 7.1.0 e versões posteriores incluem o conector, então você não precisa criar ou baixar nenhum pacote adicional. Para obter mais informações sobre o conector, consulte a [spark-sql-kinesis-connector página em GitHub](#).

Confira a seguir um exemplo de como iniciar a execução de um trabalho com a dependência do conector do Kinesis Data Streams.

```
aws emr-serverless start-job-run \  
--application-id <APPLICATION_ID> \  
--execution-role-arn <JOB_EXECUTION_ROLE> \  
--mode 'STREAMING' \  
--job-driver '{  
  "sparkSubmit": {  
    "entryPoint": "s3://<Kinesis-streaming-script>",  
    "entryPointArguments": ["s3://<DOC-EXAMPLE-BUCKET-OUTPUT>/output"],  
    "sparkSubmitParameters": "--conf spark.executor.cores=4  
      --conf spark.executor.memory=16g  
      --conf spark.driver.cores=4  
      --conf spark.driver.memory=16g  
      --conf spark.executor.instances=3  
      --jars /usr/share/aws/kinesis/spark-sql-kinesis/lib/spark-streaming-  
sql-kinesis-connector.jar"  
  }  
}'
```

Para se conectar ao Kinesis Data Streams, configure a aplicação do EMR Sem Servidor com acesso à VPC e usar um endpoint da VPC para permitir acesso privado ou um gateway NAT para obter acesso público. Para obter mais informações, consulte [Configurar o acesso à VPC](#). Também é necessário garantir que o perfil de runtime do trabalho tenha as permissões de leitura e gravação

necessárias para acessar os fluxos de dados essenciais. Para saber mais sobre como configurar um perfil de runtime de trabalho, consulte [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#). Para obter uma lista completa de todas as permissões necessárias, consulte a [spark-sql-kinesis-connector página em GitHub](#).

Conector do Apache Kafka

O conector do Apache Kafka para streaming estruturado do Spark é um conector de código aberto da comunidade Spark, disponível em um repositório Maven. Esse conector facilita que as aplicações de streaming estruturado do Spark leiam e gravem dados no Apache Kafka autogerenciado e no Amazon Managed Streaming para Apache Kafka. Para obter mais informações sobre o conector, consulte o [Structured Streaming + Kafka Integration Guide](#) na documentação do Apache Spark.

O exemplo a seguir demonstra como incluir o conector do Kafka na solicitação de execução de trabalho.

```
aws emr-serverless start-job-run \  
--application-id <APPLICATION_ID> \  
--execution-role-arn <JOB_EXECUTION_ROLE> \  
--mode 'STREAMING' \  
--job-driver '{  
  "sparkSubmit": {  
    "entryPoint": "s3://<Kafka-streaming-script>",  
    "entryPointArguments": ["s3://<DOC-EXAMPLE-BUCKET-OUTPUT>/output"],  
    "sparkSubmitParameters": "--conf spark.executor.cores=4  
      --conf spark.executor.memory=16g  
      --conf spark.driver.cores=4  
      --conf spark.driver.memory=16g  
      --conf spark.executor.instances=3  
      --packages org.apache.spark:spark-sql-  
kafka-0-10_2.12:<KAFKA_CONNECTOR_VERSION>"  
  }  
}'
```

A versão do conector do Apache Kafka depende da versão do EMR Sem Servidor e da versão correspondente do Spark. Para encontrar a versão correta do Kafka, consulte o [Structured Streaming + Kafka Integration Guide](#).

Para usar o Amazon Managed Streaming para Apache Kafka com autenticação do IAM, inclua outra dependência para permitir que o conector do Kafka se conecte ao Amazon MSK com IAM. Para obter mais informações, consulte o [aws-msk-iam-auth repositório em GitHub](#). Também é necessário

garantir que o perfil de runtime do trabalho tenha as permissões necessárias do IAM. O exemplo a seguir demonstra como usar o conector com autenticação do IAM.

```
aws emr-serverless start-job-run \
--application-id <APPLICATION_ID> \
--execution-role-arn <JOB_EXECUTION_ROLE> \
--mode 'STREAMING' \
--job-driver '{
  "sparkSubmit": {
    "entryPoint": "s3://<Kafka-streaming-script>",
    "entryPointArguments": ["s3://<DOC-EXAMPLE-BUCKET-OUTPUT>/output"],
    "sparkSubmitParameters": "--conf spark.executor.cores=4
      --conf spark.executor.memory=16g
      --conf spark.driver.cores=4
      --conf spark.driver.memory=16g
      --conf spark.executor.instances=3
      --packages org.apache.spark:spark-sql-
kafka-0-10_2.12:<KAFKA_CONNECTOR_VERSION>,software.amazon.msk:aws-msk-iam-
auth:<MSK_IAM_LIB_VERSION>"
  }
}'
```

Para usar o conector do Kafka e a biblioteca de autenticação do IAM no Amazon MSK, configure a aplicação do EMR Sem Servidor com acesso à VPC. Suas sub-redes devem ter acesso à Internet e usar um gateway NAT para acessar as dependências do Maven. Para obter mais informações, consulte [Configurar o acesso à VPC](#). As sub-redes devem ter conectividade de rede para acessar o cluster do Kafka. Isso é verdade independentemente de o cluster do Kafka ser autogerenciado ou de você usar o Amazon Managed Streaming para Apache Kafka.

Gerenciamento de logs de trabalhos de streaming

Os trabalhos de streaming oferecem suporte à alternância de logs de aplicações e eventos do Spark, bem com à compactação de logs de eventos do Spark. Isso ajuda você a gerenciar seus recursos de forma eficaz.

Alternância de logs

Os trabalhos de streaming oferecem suporte à alternância de logs de aplicações e eventos do Spark. A alternância de logs impede que trabalhos de streaming longos gerem grandes arquivos de log que podem ocupar todo o espaço disponível em disco. A alternância de logs ajuda a economizar

armazenamento em disco e evita falhas no trabalho devido ao pouco espaço em disco. Para obter mais informações, consulte [Alternar logs](#).

Compactação de logs

Os trabalhos de streaming também oferecem suporte à compactação de logs de eventos do Spark sempre que o registro em log gerenciado estiver disponível. Para obter mais detalhes sobre o registro em log gerenciado, consulte [Log com armazenamento gerenciado](#). Os trabalhos de streaming podem ser executados por um longo tempo, e a quantidade de dados de eventos pode se acumular com o tempo e aumentar significativamente o tamanho dos arquivos de log. O servidor de histórico do Spark lê e carrega esses eventos na memória para a interface do usuário da aplicação do Spark. Esse processo pode causar altas latências e custos, especialmente se os logs de eventos armazenados no Amazon S3 forem muito grandes.

A compactação de logs reduz o tamanho do log de eventos, então o servidor de histórico do Spark não precisa carregar mais de 1 GB de logs de eventos a qualquer momento. Para obter informações, consulte [Monitoring and Instrumentation](#) na documentação do Apache Spark.

Uso das configurações do Spark ao executar trabalhos do EMR Sem Servidor

Você pode executar trabalhos do Spark em uma aplicação com o parâmetro `type` definido como SPARK. Os trabalhos devem ser compatíveis com a versão do Spark compatível com a versão de lançamento do Amazon EMR. Por exemplo, quando você executa trabalhos com o Amazon EMR versão 6.6.0, seu trabalho deve ser compatível com o Apache Spark 3.2.0. Para obter informações sobre as versões da aplicação para cada versão, consulte [Versões de lançamento do Amazon EMR Sem Servidor](#).

Parâmetros do trabalho do Spark

Ao usar a [API StartJobRun](#) para executar um trabalho do Spark, especifique os parâmetros a seguir.

Parâmetros necessários

- [Perfil de runtime de trabalhos do Spark](#)
- [Parâmetro do driver de trabalhos do Spark](#)
- [Parâmetro de substituição de configuração do Spark](#)
- [Otimização dinâmica da alocação de recursos do Spark](#)

Perfil de runtime de trabalhos do Spark

Use **executionRoleArn** para especificar o ARN do perfil do IAM que a aplicação usa para executar trabalhos do Spark. Esse perfil deve conter as seguintes permissões:

- Leitura dos buckets do S3 ou de outras fontes de dados em que os dados residem
- Leia os buckets ou prefixos do S3 em que seu PySpark script ou arquivo JAR reside
- Gravação nos buckets do S3 onde você pretende gravar a saída final
- Gravação de logs em um bucket ou prefixo do S3 que `S3MonitoringConfiguration` especifica
- Acesso às chaves do KMS se você usá-las para criptografar dados no bucket do S3
- Acesso ao AWS Glue Data Catalog se você usar o SparkSQL

Se o trabalho do Spark ler ou gravar dados entre outras fontes de dados, especifique as permissões apropriadas nesse perfil do IAM. Se você não fornecer essas permissões ao perfil do IAM, o trabalho poderá falhar. Para obter mais informações, consulte [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#) e [Armazenamento de logs](#).

Parâmetro do driver de trabalhos do Spark

Use **jobDriver** para fornecer entradas ao trabalho. O parâmetro do driver do trabalho aceita somente um valor para o tipo de trabalho que você deseja executar. Para um trabalho do Spark, o valor do parâmetro é `sparkSubmit`. Você pode usar esse tipo de tarefa para executar Scala PySpark, Java e qualquer outra tarefa compatível por meio do envio do Spark. Os trabalhos do Spark têm os seguintes parâmetros:

- **sparkSubmitParameters**: esses são os parâmetros adicionais do Spark que você deseja enviar para o trabalho. Use este parâmetro para substituir as propriedades padrão do Spark, como a memória do driver ou o número de executores, como aqueles definidos nos argumentos `--conf` ou `--class`.
- **entryPointArguments**: essa é uma matriz de argumentos que você deseja transferir para o arquivo principal em JAR ou Python. Você deve realizar a leitura desses parâmetros usando seu código de `EntryPoint`. Separe cada argumento da matriz por uma vírgula.
- **entryPoint**: essa é a referência no Amazon S3 ao arquivo JAR ou Python principal que você deseja executar. Se você estiver executando um JAR Scala ou Java, especifique a classe de entrada principal nos `SparkSubmitParameters` usando o argumento `--class`.

Para obter informações adicionais, consulte [Iniciar aplicações com spark-submit](#).

Parâmetro de substituição de configuração do Spark

Use **configurationOverrides** para substituir as propriedades de configuração no nível de monitoramento e no nível de aplicação. Este parâmetro aceita um objeto JSON com os seguinte dois campos:

- **monitoringConfiguration**: use esse campo para especificar o URL do Amazon S3 (`s3MonitoringConfiguration`) no qual você deseja que o trabalho do EMR Sem Servidor armazene os logs do trabalho do Spark. Certifique-se de ter criado esse bucket com o mesmo Conta da AWS que hospeda seu aplicativo e no mesmo Região da AWS local em que seu trabalho está sendo executado.
- **applicationConfiguration**: para substituir as configurações padrão de uma aplicação, você pode fornecer um objeto de configuração nesse campo. Você pode usar uma sintaxe abreviada para fornecer a configuração ou fazer referência ao objeto de configuração em um arquivo JSON. Os objetos de configuração consistem em uma classificação, propriedades e configurações opcionais aninhadas. As propriedades consistem nas configurações que você deseja substituir neste arquivo. Você pode especificar várias classificações para diversas aplicações em um único objeto JSON.

Note

As classificações de configuração disponíveis variam de acordo com a versão específica do EMR Sem Servidor. Por exemplo, classificações para Log4j personalizado `spark-driver-log4j2` e `spark-executor-log4j2` estão disponíveis somente nas versões 6.8.0 e superiores.

Se você usar a mesma configuração em uma substituição de aplicação e nos parâmetros de envio do Spark, os parâmetros de envio do Spark terão prioridade. As configurações são classificadas em prioridade da seguinte forma, da mais alta para a mais baixa:

- Configuração que o EMR Sem Servidor fornece quando cria `SparkSession`.
- Configuração que você fornece como parte dos `sparkSubmitParameters` com o argumento `--conf`.
- A configuração que você fornece como parte da aplicação é substituída ao iniciar um trabalho.

- Configuração fornecida como parte da `runtimeConfiguration` ao criar uma aplicação.
- Configurações otimizadas escolhidas pelo Amazon EMR para a versão.
- Configurações padrão de código aberto para a aplicação.

Para obter mais informações sobre como declarar configurações no nível da aplicação e substituir configurações durante a execução do trabalho, consulte [Configuração padrão de aplicações do EMR Sem Servidor](#).

Otimização dinâmica da alocação de recursos do Spark

Use `dynamicAllocationOptimization` para otimizar o uso de recursos no EMR Sem Servidor. Definir essa propriedade como `true` na classificação de configuração do Spark indica que o EMR Sem Servidor otimizará a alocação de recursos do executor para melhor alinhar a taxa na qual o Spark solicita e cancela os executores com a taxa na qual o EMR Sem Servidor cria e libera trabalhadores. Ao fazer isso, o EMR Sem Servidor reutiliza de forma mais otimizada os trabalhadores em todos os estágios, resultando em menor custo ao executar trabalhos com vários estágios, mantendo a mesma performance.

Essa propriedade está disponível em todas as versões de lançamento do Amazon EMR.

O exemplo a seguir é de uma classificação de configuração com `dynamicAllocationOptimization`.

```
[
  {
    "Classification": "spark",
    "Properties": {
      "dynamicAllocationOptimization": "true"
    }
  }
]
```

Considere o seguinte se estiver usando a otimização de alocação dinâmica:

- Essa otimização está disponível nos trabalhos do Spark para os quais você habilitou a alocação dinâmica de recursos.
- Para obter a melhor eficiência de custos, sugerimos configurar um limite superior de ajuste de escala para os trabalhadores usando a configuração de nível de trabalho

`spark.dynamicAllocation.maxExecutors` ou a configuração de [capacidade máxima em nível de aplicação](#) com base na workload.

- Talvez você não veja uma melhora de custo em trabalhos mais simples. Por exemplo, se o trabalho for executado em um pequeno conjunto de dados ou terminar de ser executado em um estágio, talvez o Spark não precise de um número maior de executores ou de vários eventos de ajuste de escala.
- Trabalhos com uma sequência de um estágio grande, estágios menores e, em seguida, um estágio grande novamente podem sofrer regressão no runtime do trabalho. Como o EMR Sem Servidor usa os recursos com mais eficiência, isso pode levar a menos trabalhadores disponíveis para estágios maiores, ocasionando um runtime mais longo.

Propriedades do trabalho do spark

A tabela a seguir lista as propriedades opcionais do Spark e seus valores padrão que você pode substituir ao enviar um trabalho do Spark.

Propriedades opcionais do Spark e valores padrão

Chave	Description	Valor padrão
<code>spark.archives</code>	Uma lista separada por vírgulas de arquivos que o Spark extrai no diretório de trabalho de cada executor. Os tipos de arquivo compatíveis incluem <code>.jar</code> , <code>.tar.gz</code> , <code>.tgz</code> e <code>.zip</code> . Para especificar o nome do diretório a ser extraído, adicione <code>#</code> após o nome do arquivo que deseja extrair. Por exemplo, <code>.file.zip#directory</code>	NULL
<code>spark.authenticate</code>	Opção que ativa a autenticação das conexões internas do Spark.	TRUE

Chave	Description	Valor padrão
<code>spark.driver.cores</code>	O número de núcleos que o driver usa.	4
<code>spark.driver.extraJavaOptions</code>	Opções extras de Java para o driver do Spark.	NULL
<code>spark.driver.memory</code>	A quantidade de memória que o driver usa.	14G
<code>spark.dynamicAllocation.enabled</code>	Opção que ativa a alocação dinâmica de recursos. Essa opção aumenta ou reduz a escala verticalmente do número de executores registrados na aplicação com base na workload.	TRUE
<code>spark.dynamicAllocation.executorIdleTimeout</code>	O tempo que um executor pode permanecer inativo antes que o Spark o remova. Isso só se aplica se você ativar a alocação dinâmica.	60 segundos
<code>spark.dynamicAllocation.initialExecutors</code>	O número inicial de executores a serem executados se você ativar a alocação dinâmica.	3
<code>spark.dynamicAllocation.maxExecutors</code>	O limite superior do número de executores se você ativar a alocação dinâmica.	Nas versões 6.10.0 e posteriores, <i>infinity</i> Na versão 6.9.0 e inferiores, 100
<code>spark.dynamicAllocation.minExecutors</code>	O limite inferior do número de executores se você ativar a alocação dinâmica.	0

Chave	Description	Valor padrão
<code>spark.emr-serverless.allocation.batch.size</code>	O número de contêineres a serem solicitados em cada ciclo de alocação do executor. Há uma lacuna de um segundo entre cada ciclo de alocação.	20
<code>spark.emr-serverless.driver.disk</code>	O disco do driver do Spark.	20G
<code>spark.emr-serverless.driverEnv.</code> <i>[KEY]</i>	Opção que adiciona variáveis de ambiente ao driver do Spark.	NULL
<code>spark.emr-serverless.executor.disk</code>	O disco do executor do Spark.	20G
<code>spark.emr-serverless.memoryOverheadFactor</code>	Define a sobrecarga de memória a ser adicionada à memória do contêiner do driver e do executor.	0.1
<code>spark.emr-serverless.driver.disk.type</code>	O tipo de disco anexado ao driver do Spark.	Standard
<code>spark.emr-serverless.executor.disk.type</code>	O tipo de disco anexado aos executores do Spark.	Standard
<code>spark.executor.cores</code>	O número de núcleos que cada executor usa.	4
<code>spark.executor.extraJavaOptions</code>	Opções extras de Java para o executor do Spark.	NULL

Chave	Description	Valor padrão
<code>spark.executor.instances</code>	O número de contêineres do executor do Spark a serem alocados.	3
<code>spark.executor.memory</code>	A quantidade de memória que cada executor usa.	14G
<code>spark.executorEnv. [KEY]</code>	Opção que adiciona variáveis de ambiente aos executores do Spark.	NULL
<code>spark.files</code>	Uma lista separada por vírgulas de arquivos a serem colocados no diretório de trabalho de cada executor. Você pode acessar os caminhos desses arquivos no executor com <code>SparkFiles.get(fileName)</code> .	NULL
<code>spark.hadoop.hive.metastore.client.factory.class</code>	A classe de implementação da metastore do Hive.	NULL
<code>spark.jars</code>	Mais JARs para adicionar ao classpath de runtime do driver e dos executores.	NULL
<code>spark.network.crypto.enabled</code>	Opção que ativa a criptografia RPC baseada em AES. Isso inclui o protocolo de autenticação adicionado no Spark 2.2.0.	FALSE
<code>spark.sql.warehouse.dir</code>	O local padrão para bancos de dados e tabelas gerenciados.	O valor de <code>\$PWD/spark-warehouse</code>

Chave	Description	Valor padrão
<code>spark.submit.pyFiles</code>	Uma lista separada por vírgula de arquivos <code>.zip</code> , <code>.egg</code> ou <code>.py</code> para colocar no <code>PYTHONPATH</code> de aplicações do Python.	NULL

A tabela a seguir lista os parâmetros de envio padrão do Spark.

Parâmetros padrão de envio do Spark

Chave	Description	Valor padrão
<code>archives</code>	Uma lista separada por vírgulas de arquivos que o Spark extrai no diretório de trabalho de cada executor.	NULL
<code>class</code>	A classe principal da aplicação (para aplicações Java e Scala).	NULL
<code>conf</code>	Uma propriedade de configuração arbitrária do Spark.	NULL
<code>driver-cores</code>	O número de núcleos que o driver usa.	4
<code>driver-memory</code>	A quantidade de memória que o driver usa.	14G
<code>executor-cores</code>	O número de núcleos que cada executor usa.	4
<code>executor-memory</code>	A quantidade de memória que o executor usa.	14G

Chave	Description	Valor padrão
files	Uma lista separada por vírgulas de arquivos a serem colocados no diretório de trabalho de cada executor. Você pode acessar os caminhos desses arquivos no executor com <code>SparkFiles.get(fileName)</code> .	NULL
jars	Uma lista separada por vírgula de JARS a serem incluídos nos classpaths do driver e do executor.	NULL
num-executors	O número de executores a serem iniciados.	3
py-files	Uma lista separada por vírgulas de arquivos .zip, .egg ou .py para colocar no PYTHONPATH de aplicações do Python.	NULL
verbose	Opção que ativa a saída de depuração adicional.	NULL

Práticas recomendadas de configuração de recursos

Configurando recursos de driver e executor por meio da API StartJobRun

Note

Os núcleos do driver e do executor do Spark e as propriedades de memória, se especificados, devem ser especificados diretamente na solicitação da StartJobRun API.

Configurar seus recursos dessa maneira garante que o EMR Sem Servidor possa alocar os recursos corretos antes de executar o trabalho. Isso contrasta com as configurações fornecidas no script do usuário, como no arquivo .py ou .jar, que são avaliadas tarde demais, pois os operadores do driver e do executor às vezes são pré-provisionados antes do início da execução do script. Há duas formas compatíveis de configurar esses recursos durante o envio de trabalhos:

Opção 1: Usar sparkSubmitParameters

```
"jobDriver": {
  "sparkSubmit": {
    "entryPoint": "s3://your-script-path.py",
    "sparkSubmitParameters": "-conf spark.driver.memory=4g \
-CONF spark.driver.cores=2 \
-CONF spark.executor.memory=8g \
-CONF spark.executor.cores=4"
  }
}
```

Opção 2: usar configurationOverrides para a classificação spark-defaults

```
"configurationOverrides": {
  "applicationConfiguration": [
    {
      "classification": "spark-defaults",
      "properties": {
        "spark.driver.memory": "4g",
        "spark.driver.cores": "2",
        "spark.executor.memory": "8g",
        "spark.executor.cores": "4"
      }
    }
  ]
}
```

Exemplos do Spark

O exemplo a seguir mostra como usar a API `StartJobRun` para executar um script do Python. Para obter um end-to-end tutorial que usa esse exemplo, consulte [Conceitos básicos do Amazon EMR Sem Servidor](#). Você pode encontrar outros exemplos de como executar PySpark trabalhos e adicionar dependências do Python no repositório de amostras sem servidor do [EMR](#). [GitHub](#)

```
aws emr-serverless start-job-run \
  --application-id application-id \
  --execution-role-arn job-role-arn \
  --job-driver '{
    "sparkSubmit": {
      "entryPoint": "s3://us-east-1.elasticmapreduce/emr-containers/samples/
wordcount/scripts/wordcount.py",
      "entryPointArguments": ["s3://amzn-s3-demo-destination-bucket/
wordcount_output"],
      "sparkSubmitParameters": "--conf spark.executor.cores=1 --conf
spark.executor.memory=4g --conf spark.driver.cores=1 --conf spark.driver.memory=4g --
conf spark.executor.instances=1"
    }
  }'
```

O exemplo a seguir mostra como usar a API StartJobRun para executar um JAR do Spark.

```
aws emr-serverless start-job-run \
  --application-id application-id \
  --execution-role-arn job-role-arn \
  --job-driver '{
    "sparkSubmit": {
      "entryPoint": "/usr/lib/spark/examples/jars/spark-examples.jar",
      "entryPointArguments": ["1"],
      "sparkSubmitParameters": "--class org.apache.spark.examples.SparkPi --conf
spark.executor.cores=4 --conf spark.executor.memory=20g --conf spark.driver.cores=4 --
conf spark.driver.memory=8g --conf spark.executor.instances=1"
    }
  }'
```

Uso das configurações do Hive ao executar trabalhos do EMR Sem Servidor

Você pode executar trabalhos do Hive em uma aplicação com o parâmetro `type` definido como `HIVE`. Os trabalhos devem ser compatíveis com a versão do Hive compatível com a versão de lançamento do Amazon EMR. Por exemplo, quando você executa trabalhos em uma aplicação com o Amazon EMR versão 6.6.0, o trabalho deve ser compatível com o Apache Hive 3.1.2. Para obter informações sobre as versões da aplicação para cada versão, consulte [Versões de lançamento do Amazon EMR Sem Servidor](#).

Parâmetros de trabalho do Hive

Ao usar a [API StartJobRun](#) para executar um trabalho do Hive, especifique os parâmetros a seguir.

Parâmetros necessários

- [Perfil de runtime do trabalho do Hive](#)
- [Parâmetro do driver de trabalho do Hive](#)
- [Parâmetro de substituição da configuração do Hive](#)

Perfil de runtime do trabalho do Hive

Use **executionRoleArn** para especificar o ARN do perfil do IAM que a aplicação usa para executar trabalhos do Hive. Esse perfil deve conter as seguintes permissões:

- Leitura dos buckets do S3 ou de outras fontes de dados em que os dados residem
- Leitura dos buckets ou prefixos do S3 em que residem o arquivo de consulta do Hive e o arquivo de consulta inicial
- Leitura e gravação nos buckets do S3 onde residem o diretório temporário do Hive e o diretório do warehouse do Hive Metastore
- Gravação nos buckets do S3 onde você pretende gravar a saída final
- Gravação de logs em um bucket ou prefixo do S3 especificado por `S3MonitoringConfiguration`
- Acesso às chaves do KMS se você usá-las para criptografar dados no bucket do S3
- Acesso ao catálogo de dados do AWS Glue

Se o trabalho do Hive ler ou gravar dados entre outras fontes de dados, especifique as permissões apropriadas nesse perfil do IAM. Se você não fornecer essas permissões ao perfil do IAM, seu trabalho poderá falhar. Para obter mais informações, consulte [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#).

Parâmetro do driver de trabalho do Hive

Use **jobDriver** para fornecer entradas ao trabalho. O parâmetro do driver do trabalho aceita somente um valor para o tipo de trabalho que você deseja executar. Quando você especifica `hive`

como o tipo de trabalho, o EMR Sem Servidor passa uma consulta do Hive para o parâmetro `jobDriver`. Os trabalhos do Hive têm os seguintes parâmetros:

- **query**: essa é a referência no Amazon S3 ao arquivo de consulta do Hive que você deseja executar.
- **parameters**: essas são as propriedades adicionais de configuração do Hive que você deseja substituir. Para substituir propriedades, passe-as para esse parâmetro como `--hiveconf property=value`. Para substituir variáveis, passe-as para esse parâmetro como `--hivevar key=value`.
- **initQueryFile**: esse é o arquivo de consulta de inicialização do Hive. O Hive executa esse arquivo antes da consulta e pode usá-lo para inicializar tabelas.

Parâmetro de substituição da configuração do Hive

Use **configurationOverrides** para substituir as propriedades de configuração no nível de monitoramento e no nível de aplicação. Esse parâmetro aceita um objeto JSON com os seguintes dois campos:

- **monitoringConfiguration**: use esse campo para especificar o URL do Amazon S3 (`s3MonitoringConfiguration`) no qual você deseja que o trabalho do EMR Sem Servidor armazene os logs do trabalho do Hive. Certifique-se de criar esse bucket com o mesmo Conta da AWS que hospeda seu aplicativo e no mesmo Região da AWS local em que seu trabalho está sendo executado.
- **applicationConfiguration**: Você pode fornecer um objeto de configuração neste campo para substituir as configurações padrão das aplicações. Você pode usar uma sintaxe abreviada para fornecer a configuração ou fazer referência ao objeto de configuração em um arquivo JSON. Os objetos de configuração consistem em uma classificação, propriedades e configurações opcionais aninhadas. As propriedades consistem nas configurações que você deseja substituir neste arquivo. Você pode especificar várias classificações para diversas aplicações em um único objeto JSON.

Note

As classificações de configuração disponíveis variam de acordo com a versão específica do EMR Sem Servidor. Por exemplo, classificações para Log4j personalizado `spark-`

`driver-log4j2` e `spark-executor-log4j2` estão disponíveis somente nas versões 6.8.0 e superiores.

Se você transferir a mesma configuração em uma substituição de aplicação e nos parâmetros do Hive, os parâmetros do Hive terão prioridade. A lista a seguir classifica as configurações da prioridade mais alta para a mais baixa.

- Configuração fornecida como parte dos parâmetros do Hive com `--hiveconf property=value`.
- A configuração que você fornece como parte da aplicação é substituída ao iniciar um trabalho.
- Configuração fornecida como parte da `runtimeConfiguration` ao criar uma aplicação.
- Configurações otimizadas que o Amazon EMR atribui para a versão.
- Configurações padrão de código aberto da aplicação.

Para obter mais informações sobre como declarar configurações no nível da aplicação e substituir configurações durante a execução do trabalho, consulte [Configuração padrão de aplicações do EMR Sem Servidor](#).

Propriedades do trabalho do Hive

A tabela a seguir lista as propriedades obrigatórias que você configura ao enviar um trabalho do Hive.

Propriedades do trabalho do Hive

Configuração	Description
<code>hive.exec.scratchdir</code>	A localização do Amazon S3 em que o EMR Sem Servidor cria arquivos temporários durante a execução do trabalho do Hive.
<code>hive.metastore.warehouse.dir</code>	A localização dos bancos de dados do Amazon S3 para tabelas gerenciadas no Hive.

A tabela a seguir lista as propriedades opcionais do Hive e seus valores padrão que você pode substituir ao enviar um trabalho do Hive.

Propriedades opcionais do Hive e valores padrão

Configuração	Description	Valor padrão
<code>fs.s3.customAWSCredentialsProvider</code>	O provedor AWS de credenciais que você deseja usar.	<code>com.amazonaws.auth.defaultAWSCredentialsProviderChain</code>
<code>fs.s3a.aws.credentials.provider</code>	O provedor de AWS credenciais que você deseja usar com um sistema de arquivos S3A.	<code>com.amazonaws.auth.defaultAWSCredentialsProviderChain</code>
<code>hive.auto.convert.join</code>	Opção que ativa a conversão automática de junções comuns em junções de mapas, com base no tamanho do arquivo de entrada.	TRUE
<code>hive.auto.convert.join.noconditionaltask</code>	Opção que ativa a otimização quando o Hive converte uma junção comum em uma junção de mapa com base no tamanho do arquivo de entrada.	TRUE
<code>hive.auto.convert.join.noconditionaltask.size</code>	Uma junção é convertida diretamente em uma junção de mapa abaixo desse tamanho.	O valor ideal é calculado com base na memória de tarefas do Tez
<code>hive.cbo.enable</code>	Opção que ativa otimizações baseadas em custos com a estrutura Calcite.	TRUE
<code>hive.cli.tez.session.async</code>	Opção para iniciar uma sessão em segundo plano do Tez enquanto sua consulta do Hive é compilada. Quando definido como <code>false</code> , o AM	TRUE

Configuração	Description	Valor padrão
	do Tez é iniciado após a compilação da consulta do Hive.	
<code>hive.compute.query.using.stats</code>	Opção que ativa o Hive para responder a determinadas consultas com estatísticas armazenadas na metastore. Para estatísticas básicas, defina <code>hive.stats.autogather</code> como <code>TRUE</code> . Para uma coleção mais avançada de consultas, execute <code>analyze table queries</code> .	<code>TRUE</code>
<code>hive.default.fileformat</code>	O formato de arquivo padrão para instruções <code>CREATE TABLE</code> . Você pode substituir isso explicitamente se especificar <code>STORED AS [FORMAT]</code> no comando <code>CREATE TABLE</code> .	<code>TEXTFILE</code>
<code>hive.driver.cores</code>	O número de núcleos a serem usados para o processo do driver do Hive.	<code>2</code>
<code>hive.driver.disk</code>	O tamanho do disco para o driver do Hive.	<code>20G</code>
<code>hive.driver.disk.type</code>	O tipo de disco para o driver do Hive.	<code>Standard</code>
<code>hive.tez.disk.type</code>	O tamanho do disco para os trabalhadores do Tez.	<code>Standard</code>

Configuração	Description	Valor padrão
<code>hive.driver.memory</code>	A quantidade de memória a ser usada por processo de driver do Hive. A CLI do Hive e o mestre de aplicação do Tez compartilham essa memória igualmente com 20% do espaço livre.	6G
<code>hive.emr-serverless.launch.env.[KEY]</code>	Opção para definir a variável de ambiente KEY em todos os processos específicos do Hive, como o driver do Hive, o AM do Tez e o trabalho do Tez.	
<code>hive.exec.dynamic.partition</code>	Opções que ativam partições dinâmicas em DML/DDDL.	TRUE
<code>hive.exec.dynamic.partition.mode</code>	Opção que especifica se você deseja usar o modo estrito ou o modo não estrito. No modo estrito, especifique pelo menos uma partição estática para o caso de substituir acidentalmente todas as partições. No modo não estrito, todas as partições podem ser dinâmicas.	strict
<code>hive.exec.max.dynamic.partitions</code>	O número máximo de partições dinâmicas que o Hive cria no total.	1000

Configuração	Description	Valor padrão
<code>hive.exec.max.dynamic.partitions.per.node</code>	O número máximo de partições dinâmicas que o Hive cria em cada nó mapeador e redutor.	100
<code>hive.exec.orc.split.strategy</code>	Espera um dos seguintes valores: BI, ETL ou HYBRID. Essa não é uma configuração em nível de usuário. BI especifica que você deseja gastar menos tempo na geração dividida do que na execução da consulta. ETL especifica que você deseja passar mais tempo na geração dividida. HYBRID especifica uma escolha das estratégias acima com base na heurística.	HYBRID
<code>hive.exec.reducers.bytes.per.reducer</code>	O tamanho por redutor. O padrão é 256 MB. Se o tamanho da entrada for 1G, o trabalho usa 4 redutores.	256000000
<code>hive.exec.reducers.max</code>	O número máximo de redutores.	256

Configuração	Description	Valor padrão
<code>hive.exec.stagingdir</code>	O nome do diretório que armazena os arquivos temporários que o Hive cria dentro dos locais das tabelas e no local do diretório temporário especificado na propriedade <code>hive.exec.scratchdir</code> .	<code>.hive-staging</code>
<code>hive.fetch.task.conversion</code>	Espera um dos seguintes valores: NONE, MINIMAL ou MORE. O Hive pode converter consultas selecionadas em uma única tarefa de FETCH. Isso minimiza a latência.	MORE
<code>hive.groupby.position.alias</code>	Opção que faz com que o Hive use um alias de posição da coluna em instruções GROUP BY.	FALSE
<code>hive.input.format</code>	O formato de entrada padrão. Defina como <code>HiveInputFormat</code> se encontrar problemas com <code>CombineHiveInputFormat</code> .	<code>org.apache.hadoop.hive.q1.io.CombineHiveInputFormat</code>
<code>hive.log.explain.output</code>	Opção que ativa explicações sobre a saída estendida de qualquer consulta no log do Hive.	FALSE
<code>hive.log.level</code>	O nível de registro em log do Hive.	INFO

Configuração	Description	Valor padrão
<code>hive.mapred.reduce.tasks.speculative.execution</code>	Opção que ativa o lançamento o especulativo de redutores . Compatível somente com o Amazon EMR 6.10.x e versões anteriores.	TRUE
<code>hive.max-task-containers</code>	O número máximo de contêineres simultâneos. A memória do mapeador configurada é multiplicada por esse valor para determinar a memória disponível que divide o uso da computação e da preempção de tarefas.	1000
<code>hive.merge.mapfiles</code>	Opção que faz com que arquivos pequenos sejam mesclados no final de um trabalho somente de mapa.	TRUE
<code>hive.merge.size.per.task</code>	O tamanho dos arquivos mesclados no final do trabalho.	256000000
<code>hive.merge.tezfiles</code>	Opção que ativa a mesclagem de arquivos pequenos no final de um DAG do Tez.	FALSE
<code>hive.metastore.client.factory.class</code>	O nome da classe de fábrica que produz objetos que implementam a interface <code>IMetaStoreClient</code> .	<code>com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalogHiveClientFactory</code>

Configuração	Description	Valor padrão
<code>hive.metastore.glue.catalogid</code>	Se o AWS Glue Data Catalog funcionar como um metastore , mas for executado em um local Conta da AWS diferente dos trabalhos, o ID do local em Conta da AWS que os trabalhos estão sendo executados.	NULL
<code>hive.metastore.uris</code>	O URI do Thrift que o cliente da metastore usa para se conectar à metastore remota.	NULL
<code>hive.optimize.ppd</code>	Opção que ativa o pushdown do predicado.	TRUE
<code>hive.optimize.ppd.storage</code>	Opção que ativa o pushdown de predicados para manipuladores de armazenamento.	TRUE
<code>hive.orderby.position.alias</code>	Opção que faz com que o Hive use um alias de posição da coluna em instruções ORDER BY.	TRUE
<code>hive.prewarm.enabled</code>	Opção que ativa o pré-aquecimento de contêineres do Tez.	FALSE
<code>hive.prewarm.numcontainers</code>	O número de contêineres a serem pré-aquecidos para o Tez.	10

Configuração	Description	Valor padrão
<code>hive.stats.autogather</code>	Opção que faz com que o Hive colete estatísticas básicas automaticamente durante o comando <code>INSERT OVERWRITE</code> .	TRUE
<code>hive.stats.fetch.column.stats</code>	Opção que desativa a busca de estatísticas de colunas da metastore. Uma busca de estatísticas de colunas pode ser cara quando o número de colunas é alto.	FALSE
<code>hive.stats.gather.num.threads</code>	O número de threads que os comandos de análise <code>partialscan</code> e <code>noscan</code> usam para tabelas particionadas. Isso se aplica somente aos formatos de arquivo que implementam <code>StatsProvidingRecordReader</code> (como ORC).	10
<code>hive.strict.checks.cartesian.product</code>	Opções que ativam verificações estritas de junção cartesiana. Essas verificações não permitem um produto cartesiano (uma junção cruzada).	FALSE
<code>hive.strict.checks.type.safety</code>	Opção que ativa verificações de segurança de tipo estritas e desativa a comparação de <code>bigint</code> com <code>string</code> e <code>double</code> .	TRUE

Configuração	Description	Valor padrão
<code>hive.support.quote d.identifiers</code>	Espera um valor de NONE ou COLUMN. NONE implica que somente caracteres alfanuméricos e sublinhados são válidos nos identificadores. COLUMN implica que os nomes das colunas podem conter qualquer caractere.	COLUMN
<code>hive.tez.auto.reducer.parallelism</code>	Opção que ativa o recurso de paralelismo do redutor automático do Tez. O Hive ainda estima os tamanhos dos dados e define estimativas de paralelismo. O Tez coleta amostras dos tamanhos de saída dos vértices de origem e ajusta as estimativas no runtime conforme necessário.	TRUE
<code>hive.tez.container.size</code>	A quantidade de memória a ser usada por processo de tarefa do Tez.	6144
<code>hive.tez.cpu.vcores</code>	O número de núcleos a serem usados em cada tarefa do Tez.	2
<code>hive.tez.disk.size</code>	O tamanho do disco para cada contêiner de tarefas.	20G
<code>hive.tez.input.format</code>	O formato de entrada para geração de divisões no AM do Tez.	<code>org.apache.hadoop.hive.q1.io.HiveInputFormat</code>

Configuração	Description	Valor padrão
<code>hive.tez.min.partition.factor</code>	Limite inferior de redutores que o Tez especifica quando você ativa o paralelismo do redutor automático.	0.25
<code>hive.vectorized.execution.enabled</code>	Opção que ativa o modo vetorizado de execução da consulta.	TRUE
<code>hive.vectorized.execution.reduce.enabled</code>	Opção que ativa o modo vetorizado do lado reduzido da execução de uma consulta.	TRUE
<code>javax.jdo.option.ConnectionDriverName</code>	O nome da classe de driver para uma metastore JDBC.	<code>org.apache.derby.jdbc.EmbeddedDriver</code>
<code>javax.jdo.option.ConnectionPassword</code>	A senha associada a um banco de dados da metastore.	NULL
<code>javax.jdo.option.ConnectionURL</code>	A string de conexão JDBC para uma metastore JDBC.	<code>jdbc:derby;;databaseName=metastore_db;create=true</code>
<code>javax.jdo.option.ConnectionUserName</code>	O nome de usuário associado a um banco de dados da metastore.	NULL
<code>mapreduce.input.fileinputformat.split.maxsize</code>	O tamanho máximo de uma divisão durante o cálculo da divisão quando o formato de entrada é <code>org.apache.hadoop.hive.q1.io.CombineHiveInputFormat</code> . O valor de 0 indica que não há limite.	0

Configuração	Description	Valor padrão
<code>tez.am.dag.cleanup.on.completion</code>	Opção que ativa a limpeza de dados embaralhados quando o DAG é concluído.	TRUE
<code>tez.am.emr-serverless.launch.env.[KEY]</code>	Opção para definir a variável de ambiente KEY no processo de AM do Tez. Para o AM do Tez, esse valor substitui <code>hive.emr-serverless.launch.env.[KEY]</code> .	
<code>tez.am.log.level</code>	O nível de registro em log raiz que o EMR Sem Servidor passa para o primário da aplicação do Tez.	INFO
<code>tez.am.sleep.time.before.exit.millis</code>	O EMR Sem Servidor deve enviar eventos de ATS após esse período, depois da solicitação de desligamento do AM.	0
<code>tez.am.speculation.enabled</code>	Opção que causa o lançamento especulativo de tarefas mais lentas. Isso pode ajudar a reduzir a latência do trabalho quando algumas tarefas estão sendo executadas mais lentamente devido a máquinas defeituosas ou lentas. Compatível somente com o Amazon EMR 6.10.x e versões anteriores.	FALSE

Configuração	Description	Valor padrão
<code>tez.am.task.max.failed.attempts</code>	O número máximo de tentativas que podem falhar em uma tarefa específica antes que a tarefa falhe. Esse número não conta as tentativas encerradas manualmente.	3
<code>tez.am.vertex.cleanup.height</code>	Uma distância na qual, se todos os vértices dependentes estiverem completos, o AM do Tez excluirá os dados de embaralhamento de vértices. Esse recurso é desativado quando o valor é 0. O Amazon EMR 6.8.0 e versões posteriores oferecem suporte a esse recurso.	0
<code>tez.client.asynchronous-stop</code>	Opção que faz com que o EMR Sem Servidor envie eventos de ATS antes de encerrar o driver do Hive.	FALSE
<code>tez.grouping.max-size</code>	O limite superior de tamanho de uma divisão agrupada (em bytes). Esse limite evita divisões excessivamente grandes.	1073741824
<code>tez.grouping.min-size</code>	O limite inferior de tamanho de uma divisão agrupada (em bytes). Esse limite evita muitas divisões pequenas.	16777216

Configuração	Description	Valor padrão
<code>tez.runtime.io.sort.mb</code>	O tamanho do buffer flexível quando o Tez classifica a saída é classificado.	O valor ideal é calculado com base na memória de tarefas do Tez
<code>tez.runtime.unordered.output.buffer.size-mb</code>	O tamanho do buffer a ser usado se o Tez não gravar diretamente no disco.	O valor ideal é calculado com base na memória de tarefas do Tez
<code>tez.shuffle-vertex-manager.max-src-fraction</code>	A fração das tarefas de origem que devem ser concluídas antes que o EMR Sem Servidor programe todas as tarefas para o vértice atual (no caso de uma conexão ScatterGather). O número de tarefas prontas para programação no vértice atual é escalado linearmente entre <code>min-fraction</code> e <code>max-fraction</code> . Isso padroniza o valor padrão ou <code>tez.shuffle-vertex-manager.min-src-fraction</code> , o que for maior.	0.75
<code>tez.shuffle-vertex-manager.min-src-fraction</code>	A fração das tarefas de origem que devem ser concluídas antes que o EMR Sem Servidor programe tarefas para o vértice atual (no caso de uma conexão ScatterGather).	0.25

Configuração	Description	Valor padrão
<code>tez.task.emr-serverless.launch.env.[<i>KEY</i>]</code>	Opção para definir a variável de ambiente <i>KEY</i> no processo da tarefa do Tez. Para tarefas do Tez, esse valor substitui <code>hive.emr-serverless.launch.env.[<i>KEY</i>]</code> .	
<code>tez.task.log.level</code>	O nível de registro em log raiz que o EMR Sem Servidor passa para as tarefas do Tez.	INFO
<code>tez.yarn.ats.event.flush.timeout.millis</code>	A quantidade máxima de tempo que o AM deve esperar até que os eventos sejam descarregados antes do encerramento.	300000

Exemplos de trabalho do Hive

O exemplo de código a seguir mostra como executar uma consulta do Hive com a API `StartJobRun`.

```
aws emr-serverless start-job-run \
  --application-id application-id \
  --execution-role-arn job-role-arn \
  --job-driver '{
    "hive": {
      "query": "s3://amzn-s3-demo-bucket/emr-serverless-hive/query/hive-
query.q1",
      "parameters": "--hiveconf hive.log.explain.output=false"
    }
  }' \
  --configuration-overrides '{
    "applicationConfiguration": [{
      "classification": "hive-site",
      "properties": {
```

```

        "hive.exec.scratchdir": "s3://amzn-s3-demo-bucket/emr-serverless-hive/
hive/scratch",
        "hive.metastore.warehouse.dir": "s3://amzn-s3-demo-bucket/emr-
serverless-hive/hive/warehouse",
        "hive.driver.cores": "2",
        "hive.driver.memory": "4g",
        "hive.tez.container.size": "4096",
        "hive.tez.cpu.vcores": "1"
    }
}
}'

```

Você pode encontrar exemplos adicionais de como executar trabalhos do Hive no repositório [EMR GitHub Serverless Samples](#).

Resiliência de trabalhos no EMR Sem Servidor

As versões 7.1.0 e posteriores do EMR Sem Servidor incluem suporte à resiliência do trabalho, portanto, ele repete automaticamente qualquer trabalho com falha sem sua intervenção manual. Outro benefício da resiliência de trabalhos é que o EMR Sem Servidor transfere as execuções de trabalhos para uma zona de disponibilidade (AZ) diferente, caso uma AZ tenha algum problema.

Para habilitar a resiliência de um trabalho, defina a política de repetição para ele. Uma política de repetição garante que o EMR Sem Servidor reinicie automaticamente um trabalho se ele falhar em algum momento. As políticas de repetição são compatíveis com trabalhos em lote e de streaming, para que você possa personalizar a resiliência do trabalho de acordo com seu caso de uso. A tabela a seguir compara os comportamentos e as diferenças da resiliência de trabalhos em lote e de streaming.

	Trabalhos em lote	Trabalhos de streaming
Comportamento padrão do	Não executa o trabalho novamente.	Sempre tenta executar o trabalho novamente, pois a aplicação cria pontos de verificação durante a execução do trabalho.
Ponto de nova tentativa	Os trabalhos em lote não têm pontos de verificação, então	Como os trabalhos de streaming oferecem suporte

	Trabalhos em lote	Trabalhos de streaming
	o EMR Sem Servidor sempre executa novamente o trabalho desde o início.	a pontos de verificação, configure a consulta de streaming e salve o estado do runtime e o progresso em um local de ponto de verificação no Amazon S3. O EMR Sem Servidor retoma a execução do trabalho a partir do ponto de verificação. Para obter mais informações, consulte Recuperação de falhas com Checkpointing na documentação do Apache Spark.
Máximo de novas tentativas	Permite, no máximo, 10 novas tentativas.	Os trabalhos de streaming têm controle integrado de prevenção contra thrash, então a aplicação para de repetir os trabalhos se continuarem falhando após uma hora. O número padrão de novas tentativas em uma hora é cinco. Você pode configurar esse número de novas tentativas entre 1 e 10. Você não pode personalizar o número máximo de tentativas. Um valor de 1 indica que não há novas tentativas.

Quando o EMR Sem Servidor tenta executar novamente um trabalho, ele também indexa o trabalho com um número de tentativas, para que você possa acompanhar o ciclo de vida de um trabalho em todas as tentativas.

Use as operações da API do EMR Serverless ou a AWS CLI para alterar a resiliência do trabalho ou acessar informações relacionadas à resiliência do trabalho. Para obter mais informações, consulte o [Guia da API do EMR Sem Servidor](#).

Por padrão, o EMR Sem Servidor não executa novamente trabalhos em lotes. Para habilitar novas tentativas de trabalhos em lotes, configure o parâmetro `maxAttempts` ao iniciar a execução de um trabalho em lote. O parâmetro `maxAttempts` é aplicável somente a trabalhos em lotes. O padrão é 1, o que significa não executar o trabalho novamente. Os valores aceitos são de 1 a 10, inclusive.

O exemplo a seguir demonstra como especificar um número máximo de 10 tentativas ao iniciar uma execução de trabalho.

```
aws emr-serverless start-job-run
--application-id <APPLICATION_ID> \
--execution-role-arn <JOB_EXECUTION_ROLE> \
--mode 'BATCH' \
--retry-policy '{
  "maxAttempts": 10
}' \
--job-driver '{
  "sparkSubmit": {
    "entryPoint": "/usr/lib/spark/examples/jars/spark-examples-does-not-
exist.jar",
    "entryPointArguments": ["1"],
    "sparkSubmitParameters": "--class org.apache.spark.examples.SparkPi"
  }
}'
```

O EMR Sem Servidor repete indefinidamente os trabalhos de streaming se eles falharem. Para evitar thrashing devido a falhas irrecuperáveis repetidas, use `maxFailedAttemptsPerHour` para configurar o controle de prevenção contra thrash em novas tentativas de trabalhos de streaming. Esse parâmetro permite especificar o número máximo de tentativas malsucedidas permitidas até uma hora antes que o EMR Sem Servidor pare de tentar novamente. O padrão é cinco. Os valores aceitos são de 1 a 10, inclusive.

```
aws emr-serverless start-job-run
--application-id <APPLICATION_ID> \
--execution-role-arn <JOB_EXECUTION_ROLE> \
--mode 'STREAMING' \
--retry-policy '{
  "maxFailedAttemptsPerHour": 7
}'
```

```
}' \  
--job-driver '{  
  "sparkSubmit": {  
    "entryPoint": "/usr/lib/spark/examples/jars/spark-examples-does-not-  
exist.jar",  
    "entryPointArguments": ["1"],  
    "sparkSubmitParameters": "--class org.apache.spark.examples.SparkPi"  
  }  
}'
```

Você também pode usar as outras operações da API de execução de trabalhos para obter informações sobre trabalhos. Por exemplo, você pode usar o parâmetro `attempt` com a operação `GetJobRun` para obter detalhes sobre uma tentativa de trabalho específica. Se você não incluir o parâmetro `attempt`, a operação retornará informações sobre a última tentativa.

```
aws emr-serverless get-job-run \  
  --job-run-id job-run-id \  
  --application-id application-id \  
  --attempt 1
```

A operação `ListJobRunAttempts` retorna informações sobre todas as tentativas relacionadas à execução de um trabalho.

```
aws emr-serverless list-job-run-attempts \  
  --application-id application-id \  
  --job-run-id job-run-id
```

A operação `GetDashboardForJobRun` cria e retorna uma URL usada para acessar o aplicativo UIs para a execução de um trabalho. O parâmetro `attempt` permite obter um URL para uma tentativa específica. Se você não incluir o parâmetro `attempt`, a operação retornará informações sobre a última tentativa.

```
aws emr-serverless get-dashboard-for-job-run \  
  --application-id application-id \  
  --job-run-id job-run-id \  
  --attempt 1
```

Monitoramento de um trabalho com uma política de repetição

O suporte à resiliência de trabalhos também adiciona o novo evento Nova tentativa de execução de trabalho do EMR Sem Servidor. O EMR Sem servidor publica esse evento em cada nova tentativa do trabalho. Você pode usar essa notificação para rastrear novas tentativas do trabalho. Para obter mais informações sobre eventos, consulte os [EventBridge eventos da Amazon](#).

Logs com política de nova tentativa

Toda vez que o EMR Sem Servidor repete um trabalho, a tentativa gera seu próprio conjunto de logs. Para garantir que o EMR Serverless possa entregar com sucesso esses registros para o Amazon S3 e a Amazon sem CloudWatch sobrescrever nenhum, o EMR Serverless adiciona um prefixo ao formato do caminho de log do S3 e do nome do stream de log para incluir o número da tentativa do trabalho. CloudWatch

Confira a seguir um exemplo de como é o formato.

```
'/applications/<applicationId>/jobs/<jobId>/attempts/<attemptNumber>/'.
```

Esse formato garante que o EMR Serverless publique todos os registros de cada tentativa do trabalho em seu próprio local designado no Amazon S3 e. CloudWatch Para obter mais detalhes, consulte [Armazenamento de logs](#).

Note

O EMR Sem servidor usa esse formato de prefixo somente com todos os trabalhos de streaming e em lote que tenham a repetição habilitada.

Configuração da metastore para EMR Sem Servidor

Uma metastore do Hive é um local centralizado que armazena informações estruturais sobre tabelas, incluindo esquemas, nomes de partições e tipos de dados. Com o EMR Sem Servidor, você pode manter esses metadados da tabela em um metastore que tenha acesso aos seus trabalhos.

Você tem duas opções para uma metastore do Hive:

- Catálogo de dados do AWS Glue
- Uma metastore externa do Apache Hive

Usando o AWS Glue Data Catalog como metastore

Você pode configurar suas tarefas do Spark e do Hive para usar o AWS Glue Data Catalog como seu metastore. Recomendamos essa configuração quando precisar de uma metastore persistente ou de uma metastore compartilhada por diferentes aplicações, serviços, aplicações ou Contas da AWS. Para obter mais informações sobre o Catálogo de Dados, consulte [Preenchendo o Catálogo de Dados do AWS Glue](#). Para obter informações sobre os preços do AWS Glue, consulte os [preços do AWS Glue](#).

Você pode configurar sua tarefa do EMR Serverless para usar o AWS Glue Data Catalog no Conta da AWS mesmo aplicativo ou em um diferente. Conta da AWS

Configurar o catálogo de dados AWS Glue

Para configurar o Data Catalog, escolha o tipo de aplicação do EMR Sem Servidor que você deseja usar.

Spark

Quando você usa o EMR Studio para executar seus trabalhos com aplicativos EMR Serverless Spark, o AWS Glue Data Catalog é o metastore padrão.

Ao usar SDKs ou AWS CLI, defina a `spark.hadoop.hive.metastore.client.factory.class` configuração com `com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalogHiveClientFactory` nos `sparkSubmit` parâmetros da execução do seu trabalho. O exemplo a seguir mostra como configurar o Data Catalog com a AWS CLI.

```
aws emr-serverless start-job-run \  
  --application-id application-id \  
  --execution-role-arn job-role-arn \  
  --job-driver '{  
    "sparkSubmit": {  
      "entryPoint": "s3://amzn-s3-demo-bucket/code/pyspark/  
extreme_weather.py",  
      "sparkSubmitParameters": "--conf  
spark.hadoop.hive.metastore.client.factory.class=com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalogHiveClientFactory  
--conf spark.driver.cores=1 --conf spark.driver.memory=3g --conf  
spark.executor.cores=4 --conf spark.executor.memory=3g"  
    }  
  }'
```

Como alternativa, você pode definir essa configuração ao criar uma `SparkSession` no código do Spark.

```
from pyspark.sql import SparkSession

spark = (
    SparkSession.builder.appName("SparkSQL")
        .config(
            "spark.hadoop.hive.metastore.client.factory.class",
            "com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalogHiveClientFactory",
        )
        .enableHiveSupport()
        .getOrCreate()
)

# we can query tables with SparkSQL
spark.sql("SHOW TABLES").show()

# we can also them with native Spark
print(spark.catalog.listTables())
```

Hive

Para aplicações do Hive no EMR Sem Servidor, o Data Catalog é a metastore padrão. Ou seja, quando você executa trabalhos em um aplicativo EMR Serverless Hive, o Hive registra as informações do metastore no Catálogo de Dados da mesma forma que seu aplicativo. Conta da AWS Você não precisa de uma nuvem privada virtual (VPC) para usar o Data Catalog como metastore.

Para acessar as tabelas do metastore do Hive, adicione as políticas necessárias do AWS Glue descritas em [Como configurar permissões](#) do IAM para o Glue. AWS

Configure o acesso entre contas para o EMR Serverless AWS e o Glue Data Catalog

Para configurar o acesso entre contas para o EMR Serverless, primeiro faça login no seguinte: Contas da AWS

- AccountA— É Conta da AWS onde você criou um aplicativo EMR Serverless.
- AccountB— Um Conta da AWS que contém um catálogo de dados do AWS Glue que você deseja que suas execuções do EMR Serverless acessem.

1. Certifique-se de que um administrador ou outra identidade autorizada em AccountB anexe uma política de recursos ao Data Catalog em AccountB. Essa política concede permissões específicas do AccountA entre contas para realizar operações em recursos no catálogo do AccountB.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "glue:GetDatabase",
        "glue:CreateDatabase",
        "glue:GetDataBases",
        "glue:CreateTable",
        "glue:GetTable",
        "glue:UpdateTable",
        "glue>DeleteTable",
        "glue:GetTables",
        "glue:GetPartition",
        "glue:GetPartitions",
        "glue:CreatePartition",
        "glue:BatchCreatePartition",
        "glue:GetUserDefinedFunctions"
      ],
      "Resource": [
        "arn:aws:glue:*:123456789012:catalog"
      ],
      "Sid": "AllowGLUEGetdatabase"
    }
  ]
}
```

2. Adicione uma política do IAM ao perfil de runtime de trabalhos do EMR Sem Servidor em AccountA para que esse perfil possa acessar os recursos do Data Catalog em AccountB.

JSON

```
{
```

```

"Version": "2012-10-17",
"Statement": [
  {
    "Effect": "Allow",
    "Action": [
      "glue:GetDatabase",
      "glue:CreateDatabase",
      "glue:GetDataBases",
      "glue:CreateTable",
      "glue:GetTable",
      "glue:UpdateTable",
      "glue>DeleteTable",
      "glue:GetTables",
      "glue:GetPartition",
      "glue:GetPartitions",
      "glue:CreatePartition",
      "glue:BatchCreatePartition",
      "glue:GetUserDefinedFunctions"
    ],
    "Resource": [
      "arn:aws:glue:*:123456789012:catalog"
    ],
    "Sid": "AllowGLUEGetdatabase"
  }
]
}

```

3. Inicie a execução do trabalho. Essa etapa é um pouco diferente dependendo do tipo de aplicação do EMR Sem Servidor do AccountA.

Spark

Passa a propriedade `spark.hadoop.hive.metastore.glue.catalogid` no `sparkSubmitParameters` conforme mostrado no exemplo a seguir. Substitua *AccountB-catalog-id* pelo ID do Data Catalog em AccountB.

```

aws emr-serverless start-job-run \
--application-id "application-id" \
--execution-role-arn "job-role-arn" \
--job-driver '{
  "sparkSubmit": {
    "entryPoint": "s3://amzn-s3-demo-bucket/scripts/test.py",

```

```

    "sparkSubmitParameters": "--conf
spark.hadoop.hive.metastore.client.factory.class=com.amazonaws.glue.catalog.metastore.A
--conf spark.hadoop.hive.metastore.glue.catalogid=AccountB-catalog-
id --conf spark.executor.cores=1 --conf spark.executor.memory=1g
--conf spark.driver.cores=1 --conf spark.driver.memory=1g --conf
spark.executor.instances=1"
  }
}' \
--configuration-overrides '{
  "monitoringConfiguration": {
    "s3MonitoringConfiguration": {
      "logUri": "s3://amzn-s3-demo-bucket/logs/"
    }
  }
}'

```

Hive

Defina a propriedade `hive.metastore.glue.catalogid` na classificação `hive-site` conforme mostrado no exemplo a seguir. Substitua *AccountB-catalog-id* pelo ID do Data Catalog em AccountB.

```

aws emr-serverless start-job-run \
--application-id "application-id" \
--execution-role-arn "job-role-arn" \
--job-driver '{
  "hive": {
    "query": "s3://amzn-s3-demo-bucket/hive/scripts/create_table.sql",
    "parameters": "--hiveconf hive.exec.scratchdir=s3://amzn-s3-demo-bucket/
hive/scratch --hiveconf hive.metastore.warehouse.dir=s3://amzn-s3-demo-bucket/
hive/warehouse"
  }
}' \
--configuration-overrides '{
  "applicationConfiguration": [{
    "classification": "hive-site",
    "properties": {
      "hive.metastore.glue.catalogid": "AccountB-catalog-id"
    }
  }]
}'

```

Considerações ao usar o AWS Glue Data Catalog

Você pode adicionar auxiliares JARs ADD JAR em seus scripts do Hive. Para considerações adicionais, consulte Considerações ao usar o [AWS Glue Data Catalog](#).

Uso de uma metastore externa do Hive

Você pode configurar os trabalhos do Spark e do Hive no EMR Sem Servidor para se conectarem a uma metastore externa do Hive, como Amazon Aurora ou Amazon RDS para MySQL. Esta seção descreve como configurar uma metastore do Hive no Amazon RDS, uma VPC e trabalhos do EMR Sem Servidor para usar uma metastore externa.

Criação de uma metastore externa do Hive

1. Crie uma Amazon Virtual Private Cloud (Amazon VPC) com sub-redes privadas seguindo as instruções em [Crie uma VPC](#).
2. Crie sua aplicação do EMR Sem Servidor com a nova Amazon VPC e sub-redes privadas. Quando você configura a aplicação do EMR Sem Servidor com uma VPC, ela primeiro provisiona uma interface de rede elástica para cada sub-rede especificada. Em seguida, ela anexa o grupo de segurança especificado a essa interface de rede. Isso dá controle de acesso à aplicação. Para obter mais detalhes sobre como configurar a VPC, consulte [Configuração do acesso à VPC para que aplicações do EMR Sem Servidor se conectem aos dados](#).
3. Crie um banco de dados MySQL ou Aurora PostgreSQL em uma sub-rede privada na Amazon VPC. Para obter informações sobre como criar um banco de dados do Amazon RDS, consulte [Criar uma instância de banco de dados do Amazon RDS](#).
4. Modifique o grupo de segurança do banco de dados MySQL ou Aurora para permitir conexões JDBC do grupo de segurança do EMR Sem Servidor seguindo as etapas em [Modificar uma instância de banco de dados do Amazon RDS](#). Adicione uma regra para tráfego de entrada ao grupo de segurança do RDS usando um dos grupos de segurança do EMR Sem Servidor.

Tipo	Protocolo	Intervalo de portas	Fonte
Todos os TCP	TCP	3306	emr-serve rless-sec urity-group

Configurar opções do Spark

Uso do JDBC

Para configurar a aplicação do Spark no EMR Sem Servidor para se conectar a uma metastore do Hive baseada em uma instância do Amazon RDS para MySQL ou do Amazon Aurora MySQL, use uma conexão JDBC. Passe o `mariadb-connector-java.jar` com `--jars` nos parâmetros `spark-submit` da execução do trabalho.

```
aws emr-serverless start-job-run \
  --application-id "application-id" \
  --execution-role-arn "job-role-arn" \
  --job-driver '{
    "sparkSubmit": {
      "entryPoint": "s3://amzn-s3-demo-bucket/scripts/spark-jdbc.py",
      "sparkSubmitParameters": "--jars s3://amzn-s3-demo-bucket/mariadb-
connector-java.jar
      --conf
spark.hadoop.javax.jdo.option.ConnectionDriverName=org.mariadb.jdbc.Driver
      --conf spark.hadoop.javax.jdo.option.ConnectionUserName=<connection-user-
name>
      --conf spark.hadoop.javax.jdo.option.ConnectionPassword=<connection-
password>
      --conf spark.hadoop.javax.jdo.option.ConnectionURL=<JDBC-Connection-
string>
      --conf spark.driver.cores=2
      --conf spark.executor.memory=10G
      --conf spark.driver.memory=6G
      --conf spark.executor.cores=4"
    }
  }' \
  --configuration-overrides '{
    "monitoringConfiguration": {
      "s3MonitoringConfiguration": {
        "logUri": "s3://amzn-s3-demo-bucket/spark/logs/"
      }
    }
  }'
```

O exemplo de código a seguir é um script de ponto de entrada do Spark que interage com uma metastore do Hive no Amazon RDS.

```

from os.path import expanduser, join, abspath
from pyspark.sql import SparkSession
from pyspark.sql import Row
# warehouse_location points to the default location for managed databases and tables
warehouse_location = abspath('spark-warehouse')
spark = SparkSession \
    .builder \
    .config("spark.sql.warehouse.dir", warehouse_location) \
    .enableHiveSupport() \
    .getOrCreate()
spark.sql("SHOW DATABASES").show()
spark.sql("CREATE EXTERNAL TABLE `sampledb`.`sparknyctaxi`(`dispatching_base_num`
  string, `pickup_datetime` string, `dropoff_datetime` string, `pulocationid` bigint,
  `dolocationid` bigint, `sr_flag` bigint) STORED AS PARQUET LOCATION 's3://<s3 prefix>/
nyctaxi_parquet/'")
spark.sql("SELECT count(*) FROM sampledb.sparknyctaxi").show()
spark.stop()

```

Uso do serviço Thrift

Você pode configurar a aplicação do Hive no EMR Sem Servidor para se conectar a uma metastore do Hive com base em uma instância do Amazon RDS para MySQL ou do Amazon Aurora MySQL. Para fazer isso, execute um servidor do Thrift no nó primário de um cluster existente do Amazon EMR. Essa opção é ideal se você já tem um cluster do Amazon EMR com um servidor Thrift que deseja usar para simplificar as configurações de trabalho do EMR Sem Servidor.

```

aws emr-serverless start-job-run \
  --application-id "application-id" \
  --execution-role-arn "job-role-arn" \
  --job-driver '{
    "sparkSubmit": {
      "entryPoint": "s3://amzn-s3-demo-bucket/thriftscript.py",
      "sparkSubmitParameters": "--jars s3://amzn-s3-demo-bucket/mariadb-
connector-java.jar
      --conf spark.driver.cores=2
      --conf spark.executor.memory=10G
      --conf spark.driver.memory=6G
      --conf spark.executor.cores=4"
    }
  }' \
  --configuration-overrides '{
    "monitoringConfiguration": {

```

```

        "s3MonitoringConfiguration": {
            "logUri": "s3://amzn-s3-demo-bucket/spark/logs/"
        }
    }
}'

```

O exemplo de código a seguir é um script de ponto de entrada (`thriftscript.py`) que usa o protocolo Thrift para se conectar a uma metastore do Hive. Observe que a propriedade `hive.metastore.uris` precisa ser configurada para leitura em uma metastore externa do Hive.

```

from os.path import expanduser, join, abspath
from pyspark.sql import SparkSession
from pyspark.sql import Row
# warehouse_location points to the default location for managed databases and tables
warehouse_location = abspath('spark-warehouse')
spark = SparkSession \
    .builder \
    .config("spark.sql.warehouse.dir", warehouse_location) \
    .config("hive.metastore.uris", "thrift://thrift-server-host:thrift-server-port") \
    .enableHiveSupport() \
    .getOrCreate()
spark.sql("SHOW DATABASES").show()
spark.sql("CREATE EXTERNAL TABLE sampledb.`sparknyctaxi`(`dispatching_base_num`
  string, `pickup_datetime` string, `dropoff_datetime` string, `polocationid` bigint,
  `dolocationid` bigint, `sr_flag` bigint) STORED AS PARQUET LOCATION 's3://<s3 prefix>/
nyctaxi_parquet/'")
spark.sql("SELECT * FROM sampledb.sparknyctaxi").show()
spark.stop()

```

Configuração de opções do Hive

Uso do JDBC

Se você quiser especificar uma localização externa do banco de dados do Hive em uma instância do Amazon RDS MySQL ou do Amazon Aurora, você pode substituir a configuração padrão da metastore.

Note

No Hive, você pode realizar várias gravações em tabelas da metastore ao mesmo tempo. Se você compartilhar informações da metastore entre dois trabalhos, certifique-se de não gravar

na mesma tabela da metastore simultaneamente, a menos que esteja gravando em partições diferentes da mesma tabela da metastore.

Defina as configurações a seguir na classificação `hive-site` para ativar a metastore externa do Hive.

```
{
  "classification": "hive-site",
  "properties": {
    "hive.metastore.client.factory.class":
"org.apache.hadoop.hive.ql.metadata.SessionHiveMetaStoreClientFactory",
    "javax.jdo.option.ConnectionDriverName": "org.mariadb.jdbc.Driver",
    "javax.jdo.option.ConnectionURL": "jdbc:mysql://db-host:db-port/db-name",
    "javax.jdo.option.ConnectionUserName": "username",
    "javax.jdo.option.ConnectionPassword": "password"
  }
}
```

Uso de um servidor Thrift

Você pode configurar seu aplicativo EMR Serverless Hive para conectar-se a um metastore Hive baseado em um Amazon RDS for MySQL ou Amazon Aurora My. SQL Instance Para fazer isso, execute um servidor Thrift no nó principal de um cluster existente do Amazon EMR. Essa opção é ideal se você já tem um cluster do Amazon EMR que executa um servidor Thrift e deseja usar suas configurações de trabalho do EMR Sem Servidor.

Defina as configurações a seguir na classificação `hive-site` para que o EMR Sem Servidor possa acessar a metastore remota do Thrift. Observe que você deve definir a propriedade `hive.metastore.uris` para leitura em uma metastore externa do Hive.

```
{
  "classification": "hive-site",
  "properties": {
    "hive.metastore.client.factory.class":
"org.apache.hadoop.hive.ql.metadata.SessionHiveMetaStoreClientFactory",
    "hive.metastore.uris": "thrift://thrift-server-host:thrift-server-port"
  }
}
```

Trabalhando com a hierarquia de vários catálogos do AWS Glue no EMR Serverless

Você pode configurar seus aplicativos EMR Serverless para trabalhar com a hierarquia de vários catálogos do AWS Glue. O exemplo a seguir mostra como usar o EMR-S Spark com a hierarquia de vários catálogos do AWS Glue.

Para saber mais sobre a hierarquia de vários catálogos, consulte Como [trabalhar com uma hierarquia de vários catálogos no AWS Glue Data Catalog with Spark no Amazon EMR](#).

Usando o Redshift Managed Storage (RMS) com o Iceberg e o Glue Data Catalog AWS

A seguir, demonstramos como configurar o Spark para integração com um AWS Glue Data Catalog com o Iceberg:

```
aws emr-serverless start-job-run \  
  --application-id application-id \  
  --execution-role-arn job-role-arn \  
  --job-driver '{  
    "sparkSubmit": {  
      "entryPoint": "s3://amzn-s3-demo-bucket/myscript.py",  
      "sparkSubmitParameters": "--conf spark.sql.catalog.nfgac_rms =  
org.apache.iceberg.spark.SparkCatalog  
      --conf spark.sql.catalog.rms.type=glue  
      --conf spark.sql.catalog.rms.glue.id=Glue RMS catalog ID  
      --conf spark.sql.defaultCatalog=rms  
      --conf  
spark.sql.extensions=org.apache.iceberg.spark.extensions.IcebergSparkSessionExtensions"  
    }  
  }'
```

Um exemplo de consulta de uma tabela no catálogo, após a integração:

```
SELECT * FROM my_rms_schema.my_table
```

Usando o Redshift Managed Storage (RMS) com a API REST Iceberg e o Glue Data Catalog AWS

Veja a seguir como configurar o Spark para funcionar com o catálogo Iceberg REST:

```
aws emr-serverless start-job-run \  
--application-id application-id \  
--execution-role-arn job-role-arn \  
--job-driver '{  
"sparkSubmit": {  
"entryPoint": "s3://amzn-s3-demo-bucket/myscript.py",  
  "sparkSubmitParameters": "  
    --conf spark.sql.catalog.rms=org.apache.iceberg.spark.SparkCatalog  
    --conf spark.sql.catalog.rms.type=rest  
    --conf spark.sql.catalog.rms.warehouse=Glue RMS catalog ID  
    --conf spark.sql.catalog.rms.uri=Glue endpoint URI/iceberg  
    --conf spark.sql.catalog.rms.rest.sigv4-enabled=true  
    --conf spark.sql.catalog.rms.rest.signing-name=glue  
    --conf  
    spark.sql.extensions=org.apache.iceberg.spark.extensions.IcebergSparkSessionExtensions"  
  }  
}'
```

Um exemplo de consulta de uma tabela no catálogo:

```
SELECT * FROM my_rms_schema.my_table
```

Considerações sobre o uso de uma metastore externa

- Você pode configurar bancos de dados compatíveis com o JDBC do MariaDB como metastore. Exemplos desses bancos de dados são o RDS para MariaDB, MySQL e Amazon Aurora.
- As metastores não são inicializadas automaticamente. Se a metastore não for inicializada com um esquema para a versão do Hive, use a [Ferramenta de esquema do Hive](#).
- O EMR Sem Servidor não oferece suporte à autenticação do Kerberos. Você não pode usar um servidor de metastore do Thrift com autenticação do Kerberos e trabalhos do Spark ou Hive no EMR Sem Servidor.
- Você deve configurar o acesso à VPC para usar a hierarquia de vários catálogos.

Acessando dados do S3 em outra AWS conta do EMR Serverless

Você pode executar trabalhos sem servidor do Amazon EMR em uma AWS conta e configurá-los para acessar dados nos buckets do Amazon S3 que pertencem a outra conta. AWS Esta página descreve como configurar o acesso entre contas do S3 no EMR Sem Servidor.

Os trabalhos executados no EMR Serverless podem usar uma política de bucket do S3 ou uma função assumida para acessar dados no Amazon S3 a partir de uma conta diferente. AWS

Pré-requisitos

Para configurar o acesso entre contas para o Amazon EMR Serverless, conclua as tarefas enquanto estiver conectado a duas contas: AWS

- **AccountA:** essa é a conta da AWS na qual você criou uma aplicação do Amazon EMR Sem Servidor. Antes de configurar o acesso entre contas, tenha os seguintes itens prontos na conta:
 - Uma aplicação do Amazon EMR Sem Servidor em que você deseja executar trabalhos.
 - Um perfil de execução de trabalho que tem as permissões obrigatórias para executar trabalhos na aplicação. Para obter mais informações, consulte [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#).
- **AccountB:** essa é a conta da AWS que contém o bucket do S3 que você deseja que os trabalhos do Amazon EMR Sem Servidor acessem.

Uso de uma política de bucket do S3 para acessar dados entre contas do S3

Para acessar o bucket do S3 na account B usando a account A, anexe a política a seguir ao bucket do S3 na account B.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "ExamplePermissions1",
      "Effect": "Allow",
      "Action": [
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3:::my-bucket-name"
      ]
    }
  ],
}
```

```
{
  "Sid": "ExamplePermissions2",
  "Effect": "Allow",
  "Action": [
    "s3:PutObject",
    "s3:GetObject",
    "s3:DeleteObject"
  ],
  "Resource": [
    "arn:aws:s3:::my-bucket-name/*"
  ]
}
```

Para obter mais informações sobre o acesso entre contas do S3 com políticas de bucket do S3, consulte [Exemplo 2: Proprietário do bucket concedendo permissões de bucket entre contas](#) no Guia do usuário do Amazon Simple Storage Service.

Uso de um perfil assumido para acessar dados entre contas do S3

Outra forma de configurar o acesso entre contas para o Amazon EMR Serverless é com `AssumeRole` a ação do (). AWS Security Token Service AWS STS AWS STS é um serviço web global que permite solicitar credenciais temporárias com privilégios limitados para os usuários. Você pode fazer chamadas de API para o EMR Sem Servidor e o Amazon S3 com as credenciais de segurança temporárias criadas com `AssumeRole`.

As seguintes etapas ilustram como usar um perfil assumido para acessar dados entre contas do S3 no EMR Sem Servidor:

1. Crie um bucket do Amazon S3, *cross-account-bucket*, na AccountB. Para mais informações, consulte [Criar um bucket](#) no Guia do usuário do Amazon Simple Storage Service. Se quiser ter acesso entre contas ao DynamoDB, você também pode criar uma tabela do DynamoDB na AccountB. Para obter mais informações, consulte [Criar uma tabela no DynamoDB](#), no Guia do desenvolvedor do Amazon DynamoDB.
2. Crie um perfil do IAM `Cross-Account-Role-B` na AccountB que possa acessar o *cross-account-bucket*.
 - a. Faça login no Console de gerenciamento da AWS e abra o console do IAM em <https://console.aws.amazon.com/iam/>.

- b. Escolha Perfis e crie um novo perfil: `Cross-Account-Role-B`. Para obter mais informações sobre como criar perfis do IAM, consulte [Criar um perfil do IAM](#), no Guia do usuário do IAM.
- c. Crie uma política do IAM que especifique as permissões para `Cross-Account-Role-B` acessar o bucket `cross-account-bucket` do S3, como demonstra a instrução de política a seguir. Em seguida, anexe a política do IAM ao `Cross-Account-Role-B`. Para obter mais informações, consulte [Criar políticas do IAM](#) no Guia do usuário do IAM.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:*"
      ],
      "Resource": [
        "arn:aws:s3:::cross-account-bucket",
        "arn:aws:s3:::cross-account-bucket/*"
      ],
      "Sid": "AllowS3"
    }
  ]
}
```

Se o acesso ao DynamoDB for necessário, crie uma política do IAM que especifique as permissões de acesso à tabela do DynamoDB entre contas. Em seguida, anexe a política do IAM ao `Cross-Account-Role-B`. Para obter mais informações, consulte [Amazon DynamoDB: permite acesso a uma tabela específica](#) no Guia do usuário do IAM.

A seguir, é apresentada uma política para permitir acesso à tabela `CrossAccountTable` do DynamoDB.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
```

```

    "Action": [
      "dynamodb:*"
    ],
    "Resource": [
      "arn:aws:dynamodb:*:123456789012:table/CrossAccountTable"
    ],
    "Sid": "AllowDYNAMODB"
  }
]
}

```

3. Edite a relação de confiança para o perfil Cross-Account-Role-B.

- Para configurar o relacionamento de confiança para o perfil, escolha a guia Relacionamentos de confiança no console do IAM para o perfil Cross-Account-Role-B criado na Etapa 2.
- Selecione Editar relação de confiança.
- Insira o documento de política a seguir. Isso permite que Job-Execution-Role-A na AccountA assuma o perfil Cross-Account-Role-B.

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "sts:AssumeRole"
      ],
      "Resource": "arn:aws:iam::123456789012:role/Job-Execution-Role-A",
      "Sid": "AllowSTSAssumerole"
    }
  ]
}

```

4. Job-Execution-Role-A Conceda AccountA a AWS STS AssumeRole permissão para assumir Cross-Account-Role-B.

- No console do IAM da AWS conta AccountA, selecione Job-Execution-Role-A.
- Adicione a instrução de política a seguir ao Job-Execution-Role-A para permitir a ação AssumeRole no perfil Cross-Account-Role-B.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "sts:AssumeRole"
      ],
      "Resource": [
        "arn:aws:iam::123456789012:role/Cross-Account-Role-B"
      ],
      "Sid": "AllowSTSAssumerole"
    }
  ]
}
```

Exemplos de perfis assumidos

Use um único perfil assumido para acessar todos os recursos do S3 em uma conta ou, com o Amazon EMR 6.11 e superior, configurar vários perfis do IAM a serem assumidos ao acessar diferentes buckets do S3 entre contas.

Tópicos

- [Acesso aos recursos do S3 com um perfil assumido](#)
- [Acesso a recursos do S3 com vários perfis assumidos](#)

Acesso aos recursos do S3 com um perfil assumido

Note

Quando você configura um trabalho para usar um único perfil assumido, todos os recursos do S3 em todo o trabalho usam esse perfil, incluindo o script `entryPoint`.

Se você quiser usar um único perfil assumido para acessar todos os recursos do S3 na conta B, especifique as seguintes configurações:

1. Especifique a configuração `fs.s3.customAWSCredentialsProvider` do EMRFS para `com.amazonaws.emr.AssumeRoleAWSCredentialsProvider`.
2. No Spark, use `spark.emr-serverless.driverEnv.ASSUME_ROLE_CREDENTIALS_ROLE_ARN` e `spark.executorEnv.ASSUME_ROLE_CREDENTIALS_ROLE_ARN` para especificar as variáveis de ambiente no driver e nos executores.
3. No Hive, use `hive.emr-serverless.launch.env.ASSUME_ROLE_CREDENTIALS_ROLE_ARN`, `tez.am.emr-serverless.launch.env.ASSUME_ROLE_CREDENTIALS_ROLE_ARN` e `tez.task.emr-serverless.launch.env.ASSUME_ROLE_CREDENTIALS_ROLE_ARN` para especificar as variáveis de ambiente no driver do Hive, no primário de aplicação do Tez e nos contêineres de tarefas do Tez.

Os exemplos a seguir mostram como usar um perfil assumido para iniciar uma execução de trabalho do EMR Sem Servidor com acesso entre contas.

Spark

O exemplo a seguir mostra como usar um perfil assumido para iniciar a execução de um trabalho do Spark no EMR Sem Servidor com acesso entre contas ao S3.

```
aws emr-serverless start-job-run \
  --application-id application-id \
  --execution-role-arn job-role-arn \
  --job-driver '{
    "sparkSubmit": {
      "entryPoint": "entrypoint_location",
      "entryPointArguments": [":argument_1:", ":argument_2:"],
      "sparkSubmitParameters": "--conf spark.executor.cores=4 --conf
spark.executor.memory=20g --conf spark.driver.cores=4 --conf spark.driver.memory=8g
--conf spark.executor.instances=1"
    }
  }' \
  --configuration-overrides '{
    "applicationConfiguration": [{
      "classification": "spark-defaults",
      "properties": {
        "spark.hadoop.fs.s3.customAWSCredentialsProvider":
"com.amazonaws.emr.AssumeRoleAWSCredentialsProvider",
```

```

        "spark.emr-serverless.driverEnv.ASSUME_ROLE_CREDENTIALS_ROLE_ARN":
"arn:aws:iam::AccountB:role/Cross-Account-Role-B",
        "spark.executorEnv.ASSUME_ROLE_CREDENTIALS_ROLE_ARN":
"arn:aws:iam::AccountB:role/Cross-Account-Role-B"
    }
  ]]
}'

```

Hive

O exemplo a seguir mostra como usar um perfil assumido para iniciar uma execução de trabalho do Hive no EMR Sem Servidor com acesso entre contas ao S3.

```

aws emr-serverless start-job-run \
  --application-id application-id \
  --execution-role-arn job-role-arn \
  --job-driver '{
    "hive": {
      "query": "query_location",
      "parameters": "hive_parameters"
    }
  }' \
  --configuration-overrides '{
    "applicationConfiguration": [{
      "classification": "hive-site",
      "properties": {
        "fs.s3.customAWSCredentialsProvider":
"com.amazonaws.emr.serverless.credentialsprovider.AssumeRoleAWSCredentialsProvider",
        "hive.emr-serverless.launch.env.ASSUME_ROLE_CREDENTIALS_ROLE_ARN":
"arn:aws:iam::AccountB:role/Cross-Account-Role-B",
        "tez.am.emr-serverless.launch.env.ASSUME_ROLE_CREDENTIALS_ROLE_ARN":
"arn:aws:iam::AccountB:role/Cross-Account-Role-B",
        "tez.task.emr-
serverless.launch.env.ASSUME_ROLE_CREDENTIALS_ROLE_ARN":
"arn:aws:iam::AccountB:role/Cross-Account-Role-B"
      }
    }
  }'

```

Acesso a recursos do S3 com vários perfis assumidos

Com o EMR Sem Servidor nas versões 6.11.0 e posteriores, configure vários perfis do IAM a serem assumidos ao acessar diferentes buckets entre contas. Se você quiser acessar diferentes recursos do S3 com diferentes perfis assumidos na conta B, use as seguintes configurações ao iniciar a execução do trabalho:

1. Especifique a configuração `fs.s3.customAWSCredentialsProvider` do EMRFS para `com.amazonaws.emr.serverless.credentialsprovider.BucketLevelAssumeRoleCredentialsProvider`.
2. Especifique a configuração `fs.s3.bucketLevelAssumeRoleMapping` do EMRFS para definir o mapeamento do nome do bucket do S3 para o perfil do IAM na conta B a ser assumido. O valor deve estar no formato de `bucket1->role1;bucket2->role2`.

Por exemplo, use `arn:aws:iam::AccountB:role/Cross-Account-Role-B-1` para acessar o bucket `bucket1` e `arn:aws:iam::AccountB:role/Cross-Account-Role-B-2` para acessar o bucket `bucket2`. Os exemplos a seguir mostram como iniciar uma execução de trabalho do EMR Sem Servidor com acesso entre contas por meio de vários perfis assumidos.

Spark

O exemplo a seguir mostra como usar vários perfis assumidos para criar uma execução de trabalho do Spark no EMR Sem Servidor.

```
aws emr-serverless start-job-run \
  --application-id application-id \
  --execution-role-arn job-role-arn \
  --job-driver '{
    "sparkSubmit": {
      "entryPoint": "entrypoint_location",
      "entryPointArguments": [":argument_1:", ":argument_2:"],
      "sparkSubmitParameters": "--conf spark.executor.cores=4 --conf
spark.executor.memory=20g --conf spark.driver.cores=4 --conf spark.driver.memory=8g
--conf spark.executor.instances=1"
    }
  }' \
  --configuration-overrides '{
    "applicationConfiguration": [{
      "classification": "spark-defaults",
      "properties": {
```

```

        "spark.hadoop.fs.s3.customAWSCredentialsProvider":
"com.amazonaws.emr.serverless.credentialsprovider.BucketLevelAssumeRoleCredentialsProvider"
        "spark.hadoop.fs.s3.bucketLevelAssumeRoleMapping":
"bucket1->arn:aws:iam::AccountB:role/Cross-Account-Role-B-1;bucket2-
>arn:aws:iam::AccountB:role/Cross-Account-Role-B-2"
    }
  ]]
}'

```

Hive

Os exemplos a seguir mostram como usar vários perfis assumidos para criar uma execução de trabalho do Hive no EMR Sem Servidor.

```

aws emr-serverless start-job-run \
  --application-id application-id \
  --execution-role-arn job-role-arn \
  --job-driver '{
    "hive": {
      "query": "query_location",
      "parameters": "hive_parameters"
    }
  }' \
  --configuration-overrides '{
    "applicationConfiguration": [{
      "classification": "hive-site",
      "properties": {
        "fs.s3.customAWSCredentialsProvider":
"com.amazonaws.emr.serverless.credentialsprovider.AssumeRoleAWSCredentialsProvider",
        "fs.s3.bucketLevelAssumeRoleMapping": "bucket1-
>arn:aws:iam::AccountB:role/Cross-Account-Role-B-1;bucket2-
>arn:aws:iam::AccountB:role/Cross-Account-Role-B-2"
      }
    }
  }'

```

Solução de problemas no EMR Sem Servidor

Use as informações a seguir para ajudar a diagnosticar e corrigir problemas comuns que ocorrem ao trabalhar com o Amazon EMR Sem Servidor.

Tópicos

- [Erro: o trabalho falhou porque a conta atingiu o limite do serviço para o número máximo de vCPUs que pode usar simultaneamente.](#)
- [Erro: o trabalho falhou porque a aplicação excedeu as configurações de capacidade máxima.](#)
- [Erro: o trabalho falhou porque o operador não pôde ser alocado porque a aplicação excedeu a capacidade máxima.](#)
- [Erro: acesso ao S3 negado. Verifique as permissões de acesso ao S3 do perfil de runtime do trabalho nos recursos necessários do S3.](#)
- [Erro ModuleNotFoundError: Nenhum módulo nomeado<module>. Consulte o guia do usuário sobre como usar bibliotecas Python com o EMR Sem Servidor.](#)
- [Erro: Não foi possível assumir o perfil de execução <role name> porque ele não existe ou não está configurado com o relacionamento de confiança necessário.](#)

Erro: o trabalho falhou porque a conta atingiu o limite do serviço para o número máximo de vCPUs que pode usar simultaneamente.

Este erro indica que o EMR Sem Servidor não conseguiu enviar o trabalho porque a conta excedeu a capacidade máxima. Aumente a capacidade máxima da conta. Verifique seus limites de serviço nas [cotas de serviço do EMR Sem Servidor](#).

Erro: o trabalho falhou porque a aplicação excedeu as configurações de capacidade máxima.

Esse erro indica que o EMR Sem Servidor não conseguiu enviar o trabalho, pois a aplicação excedeu a capacidade máxima configurada. Aumente a capacidade máxima da aplicação.

Erro: o trabalho falhou porque o operador não pôde ser alocado porque a aplicação excedeu a capacidade máxima.

Esse erro indica que a tarefa não foi concluída. Operadores não puderam ser alocados porque a aplicação excedeu as configurações de capacidade máxima.

Erro: acesso ao S3 negado. Verifique as permissões de acesso ao S3 do perfil de runtime do trabalho nos recursos necessários do S3.

Esse erro indica que o trabalho não tem acesso a recursos do S3. Verifique se o perfil de runtime do trabalho tem permissão para acessar os recursos do S3 que o trabalho precisa usar. Para saber mais sobre perfis de runtime, consulte [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#).

Erro ModuleNotFoundError: Nenhum módulo nomeado<module>. Consulte o guia do usuário sobre como usar bibliotecas Python com o EMR Sem Servidor.

Esse erro indica que um módulo Python não estava disponível para o trabalho do Spark. Verifique se as bibliotecas Python dependentes estão disponíveis para o trabalho. Para obter mais informações sobre como empacotar bibliotecas Python, consulte [Uso de bibliotecas Python com o EMR Sem Servidor](#).

Erro: Não foi possível assumir o perfil de execução <role name> porque ele não existe ou não está configurado com o relacionamento de confiança necessário.

Esse erro indica que o perfil de runtime do trabalho que você especificou para o trabalho não existe ou que o perfil não tem um relacionamento de confiança para as permissões do EMR Sem Servidor. Para verificar se o perfil do IAM existe e validar se você configurou a política de confiança do perfil corretamente, consulte as instruções em [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#).

Habilitando a alocação de custos no nível de Job

A alocação de custos no nível do trabalho permite a atribuição granular do faturamento do EMR Serverless no nível de execução individual do trabalho, em vez de agregar todos os custos no nível do aplicativo. Quando ativado, você pode filtrar e rastrear AWS custos no Cost Explorer e nos Relatórios de Custo e Uso por execução de tarefas específicas IDs e tags associadas às execuções de tarefas, fornecendo melhor visibilidade das cobranças pelas execuções de tarefas enviadas.

Comportamento padrão do

A alocação de custos em nível de trabalho não está ativada por padrão.

Como ativar ou desativar o recurso

Você pode configurar a alocação de custos em nível de trabalho durante a criação do aplicativo ou atualizá-la para aplicativos existentes.

Especifique o `jobLevelCostAllocation` parâmetro ao criar um novo aplicativo:

```
# Enable job-level cost allocation:
aws emr-serverless create-application \
  --name "my-application" \
  --release-label "emr-7.12.0" \
  --type "SPARK" \
  --job-level-cost-allocation-configuration '{
    "enabled": true
  }'

# Disable job-level cost allocation:
aws emr-serverless create-application \
  --name "my-application" \
  --release-label "emr-7.12.0" \
  --type "SPARK" \
  --job-level-cost-allocation-configuration '{
    "enabled": false
  }'
```

Atualize o `jobLevelCostAllocationConfiguration` parâmetro para um aplicativo existente:

```
# Enable job-level cost allocation:
aws emr-serverless update-application \
  --application-id <application-id> \
  --job-level-cost-allocation-configuration '{
    "enabled": true
  }'

# Disable job-level cost allocation:
aws emr-serverless update-application \
  --application-id <application-id> \
  --job-level-cost-allocation-configuration '{
```

```
"enabled": false  
'
```

Condições e limitações

- A ativação da alocação de custos em nível de trabalho não atribui retroativamente custos às execuções de tarefas concluídas antes da ativação do recurso. As execuções de trabalhos iniciadas após a ativação do recurso terão uma atribuição granular de custo.
- O parâmetro de alocação de custos no nível do trabalho só pode ser atualizado quando um aplicativo está no estado CRIADO ou PARADO.
- Quando a alocação de custos em nível de trabalho é ativada, os custos são atribuídos às execuções de tarefas individuais, e não ao aplicativo. Para visualizar os custos agregados no nível do aplicativo, você deve aplicar tags consistentes (como application-name ou application-id) a todas as execuções de trabalhos dentro desse aplicativo e filtrar por essas tags no Cost Explorer ou nos Relatórios de Custos e Uso.

Execução de workloads interativas com o EMR Sem Servidor por meio do EMR Studio

Com as aplicações interativas do EMR Sem Servidor, execute workloads interativas para o Spark com o EMR Sem Servidor usando cadernos hospedados no EMR Studio.

Visão geral do

Uma aplicação interativa é uma aplicação do EMR Sem Servidor que tem recursos interativos habilitados. Com as aplicações interativas do Amazon EMR Sem Servidor, você pode executar workloads interativas com cadernos Jupyter gerenciados no Amazon EMR Studio. Essa ação ajuda engenheiros de dados, cientistas de dados e analistas de dados a usar o EMR Studio para executar analytics interativas com conjuntos de dados em armazenamentos de dados como o Amazon S3 e o Amazon DynamoDB.

Os casos de uso de aplicações interativas no EMR Sem Servidor incluem:

- Os engenheiros de dados usam a experiência de IDE no EMR Studio para criar um script ETL. O script ingere dados on-premises, os transforma para análise e os armazena no Amazon S3.
- Os cientistas de dados usam cadernos para explorar conjuntos de dados e treinar modelos de machine learning (ML) para detectar anomalias nos conjuntos de dados.
- Os analistas de dados exploram conjuntos de dados e criam scripts que geram relatórios diários para atualizar aplicações como painéis de negócios.

Pré-requisitos

Para usar workloads interativas com o EMR Sem Servidor, atenda aos seguintes requisitos:

- As aplicações interativas do EMR Sem Servidor são compatíveis com o Amazon EMR 6.14.0 e versões superiores.
- Para acessar sua aplicação interativa, executar as workloads enviadas e executar cadernos interativos do EMR Studio, você precisa de permissões e perfis específicos. Para obter mais informações, consulte [Permissões necessárias para workloads interativas](#).

Permissões necessárias para workloads interativas

Além das [permissões básicas necessárias para acessar o EMR Sem Servidor](#), configure permissões adicionais para a identidade ou o perfil do IAM:

Para acessar a aplicação interativa

Configure as permissões de usuário e Workspace do EMR Studio. Para obter mais informações, consulte [Configurar permissões de usuário do EMR Studio](#) no Guia de gerenciamento do Amazon EMR.

Para executar as workloads enviadas com o EMR Sem Servidor

Configuração de um perfil de runtime do trabalho. Para obter mais informações, consulte [Criação de um perfil de runtime de trabalhos](#).

Para executar os cadernos interativos do EMR Studio

Adicione as seguintes permissões adicionais à política do IAM para os usuários do Studio:

- **emr-serverless:AccessInteractiveEndpoints**: concede permissão para acessar e se conectar à aplicação interativa que você especifica como Resource. Essa permissão é necessária para se anexar a uma aplicação do EMR Sem Servidor de um Workspace do EMR Studio.
- **iam:PassRole**: concede permissão para acessar o perfil de execução do IAM que você planeja usar ao se conectar a uma aplicação. É necessária a permissão PassRole apropriada para se anexar a uma aplicação do EMR Sem Servidor de um Workspace do EMR Studio.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "EMRServerlessInteractiveAccess",
      "Effect": "Allow",
      "Action": [
        "emr-serverless:AccessInteractiveEndpoints"
      ],
      "Resource": [
        "arn:aws:emr-serverless:*:123456789012:/applications/*"
      ]
    }
  ],
}
```

```
{
  "Sid": "EMRServerlessRuntimeRoleAccess",
  "Effect": "Allow",
  "Action": [
    "iam:PassRole"
  ],
  "Resource": [
    "arn:aws:iam::123456789012:role/EMRServerlessInteractiveRole"
  ],
  "Condition": {
    "StringLike": {
      "iam:PassedToService": "emr-serverless.amazonaws.com"
    }
  }
}
```

Configuração de aplicações interativas

Use as etapas de alto nível a seguir para criar uma aplicação do EMR Sem Servidor com recursos interativos do Amazon EMR Studio no Console de gerenciamento da AWS.

1. Siga as etapas em [Conceitos básicos do Amazon EMR Sem Servidor](#) para criar uma aplicação.
2. Em seguida, inicie um Workspace do EMR Studio e se anexe a uma aplicação do EMR Sem Servidor como opção de computação. Para obter mais informações, consulte a guia Workload interativa na Etapa 2 da documentação de [Conceitos básicos sobre o EMR Sem Servidor](#).

Quando você anexa uma aplicação a um Workspace do Studio, o início da aplicação é acionado automaticamente se ainda não estiver em execução. Você também pode pré-iniciar a aplicação e mantê-la pronta antes de anexá-la ao Workspace.

Considerações sobre aplicações interativas

- As aplicações interativas do EMR Sem Servidor são compatíveis com o Amazon EMR 6.14.0 e versões superiores.

- O EMR Studio é o único cliente integrado às aplicações interativas do EMR Sem Servidor. Os seguintes recursos do EMR Studio não são compatíveis com as aplicações interativas do EMR Sem Servidor: colaboração no Workspace, SQL Explorer e execução programática de cadernos.
- Aplicações interativas são compatíveis apenas com o mecanismo Spark.
- Os aplicativos interativos oferecem suporte aos kernels Python 3 PySpark e Spark Scala.
- Você pode executar até 25 cadernos simultâneos em uma única aplicação interativa.
- Não há um endpoint ou interface de API que ofereça suporte a cadernos Jupyter auto-hospedados com aplicações interativas.
- Para uma experiência de inicialização otimizada, sugerimos configurar a capacidade pré-inicializada de drivers e executores e pré-iniciar a aplicação. Ao pré-iniciar a aplicação, você garante que ela esteja pronta quando quiser anexá-la ao Workspace.

```
aws emr-serverless start-application \  
--application-id your-application-id
```

- Por padrão, autoStopConfig está habilitado para aplicações. Isso desliga a aplicação após 30 minutos de tempo ocioso. Você pode alterar essa configuração como parte da solicitação create-application ou update-application.
- Ao usar uma aplicação interativa, sugerimos configurar uma capacidade pré-inicializada de kernels, drivers e executores para executar os cadernos. Cada sessão interativa do Spark requer um kernel e um driver, então o EMR Sem Servidor mantém um trabalhador do kernel pré-inicializado para cada driver pré-inicializado. Por padrão, o EMR Sem Servidor mantém uma capacidade pré-inicializada de um trabalhador do kernel em toda a aplicação, mesmo que você não especifique nenhuma capacidade pré-inicializada para drivers. Cada trabalhador do kernel usa 4 vCPUs e 16 GB de memória. Para obter informações atuais sobre preço, consulte a página [Preços do Amazon EMR](#).
- Você deve ter uma cota de serviço de vCPU suficiente Conta da AWS para executar cargas de trabalho interativas. Se você não executa workloads habilitadas para Lake Formation, sugerimos pelo menos 24 vCPUs. Caso execute, sugerimos no mínimo 28 vCPUs.
- O EMR Sem Servidor encerra automaticamente os kernels dos cadernos se eles ficarem inativos por mais de 60 minutos. O EMR Sem Servidor calcula o tempo ocioso do kernel a partir da última atividade concluída durante a sessão do caderno. No momento, você não pode modificar a configuração do tempo limite de inatividade do kernel.
- Para habilitar o Lake Formation com workloads interativas, defina a configuração spark.emr-serverless.lakeformation.enabled para true na classificação spark-defaults no

objeto `runtime-configuration` ao [criar uma aplicação do EMR Sem Servidor](#). Para saber mais, consulte [Habilitar o Lake Formation no Amazon EMR](#).

Execução de workloads interativas com o EMR Servidor por meio de um endpoint do Apache Livy

Com as versões 6.14.0 e superiores do Amazon EMR, você pode criar e habilitar um endpoint do Apache Livy enquanto cria uma aplicação do EMR Sem Servidor e executar workloads interativas por meio de cadernos auto-hospedados ou com um cliente personalizado. Um endpoint do Apache Livy oferece os seguintes benefícios:

- Você pode se conectar com segurança a um endpoint do Apache Livy por meio de cadernos Jupyter e gerenciar workloads do Apache Spark com a interface REST do Apache Livy.
- Use as operações da API REST do Apache Livy para aplicações Web interativas que utilizam dados das workloads do Apache Spark.

Pré-requisitos

Para usar um endpoint do Apache Livy com o EMR Sem Servidor, atenda aos seguintes requisitos:

- Conclua as etapas em [Getting started with Amazon EMR Sem Servidor](#).
- Para executar workloads interativas por meio dos endpoints do Apache Livy, você precisa de determinadas permissões e perfis. Para obter mais informações, consulte [Permissões necessárias para workloads interativas](#).

Permissões obrigatórias

Além das permissões necessárias para acessar o EMR Sem Servidor, adicione também as seguintes permissões ao perfil do IAM para acessar um endpoint do Apache Livy e executar aplicações:

- `emr-serverless:AccessLivyEndpoints`: concede permissão para acessar e se conectar à aplicação habilitada para Livy especificada como `Resource`. Você precisa dessa permissão para executar as operações da API REST disponíveis no endpoint do Apache Livy.
- `iam:PassRole`: concede permissão para acessar o perfil de execução do IAM ao criar a sessão do Apache Livy. O EMR Sem Servidor usará esse perfil para executar suas workloads.

- `emr-serverless:GetDashboardForJobRun`: concede permissão para gerar a interface do usuário do Spark Live e os links de log do driver e fornece acesso aos logs como parte dos resultados da sessão do Apache Livy.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "EMRServerlessInteractiveAccess",
      "Effect": "Allow",
      "Action": [
        "emr-serverless:AccessLivyEndpoints"
      ],
      "Resource": [
        "arn:aws:emr-serverless:*:123456789012:/applications/*"
      ]
    },
    {
      "Sid": "EMRServerlessRuntimeRoleAccess",
      "Effect": "Allow",
      "Action": [
        "iam:PassRole"
      ],
      "Resource": [
        "arn:aws:iam::123456789012:role/EMRServerlessExecutionRole"
      ],
      "Condition": {
        "StringLike": {
          "iam:PassedToService": "emr-serverless.amazonaws.com"
        }
      }
    },
    {
      "Sid": "EMRServerlessDashboardAccess",
      "Effect": "Allow",
      "Action": [
        "emr-serverless:GetDashboardForJobRun"
      ],
      "Resource": [
        "arn:aws:emr-serverless:*:123456789012:/applications/*"
      ]
    }
  ]
}
```

```
    ]
  }
]
}
```

Introdução

Para criar uma aplicação habilitada para Apache Livy e executá-la, siga estas etapas.

1. Para criar uma aplicação habilitada para Apache Livy, execute o comando a seguir.

```
aws emr-serverless create-application \  
--name my-application-name \  
--type 'application-type' \  
--release-label <Amazon EMR-release-version> \  
--interactive-configuration '{"livyEndpointEnabled": true}'
```

2. Depois que o EMR Sem Servidor criar a aplicação, inicie-a para disponibilizar o endpoint do Apache Livy.

```
aws emr-serverless start-application \  
--application-id application-id
```

Use o comando a seguir para verificar o status da aplicação. Quando o status se tornar STARTED, acesse o endpoint do Apache Livy.

```
aws emr-serverless get-application \  
--region <AWS_REGION> --application-id >application_id>
```

3. Use o seguinte URL para acessar o endpoint:

```
https://_application-id_.livy.emr-serverless-  
services._<AWS_REGION>_.amazonaws.com
```

Quando o endpoint estiver pronto, envie workloads com base no caso de uso. Você deve assinar todas as solicitações no endpoint com [o SIGv4 protocolo](#) e passar um cabeçalho de autorização. É possível usar os seguintes métodos para executar workloads:

- Cliente HTTP: envie suas operações da API de endpoint do Apache Livy com um cliente HTTP personalizado.
- Kernel do Sparkmagic: execute localmente o kernel do Sparkmagic e envie consultas interativas com os cadernos Jupyter.

Cientes HTTP

Para criar uma sessão do Apache Livy, envie `emr-serverless.session.executionRoleArn` no parâmetro `conf` do corpo da solicitação. Confira a seguir um exemplo de solicitação `POST /sessions`.

```
{
  "kind": "pyspark",
  "heartbeatTimeoutInSeconds": 60,
  "conf": {
    "emr-serverless.session.executionRoleArn": "<executionRoleArn>"
  }
}
```

A tabela a seguir descreve todas as operações disponíveis da API do Apache Livy.

Operação de API	Description
GET /sessions	Retorna uma lista de todas as sessões interativas ativas.
POST /sessions	Cria uma sessão interativa por meio do Spark ou PySpark.
OBTENHA /sessions/ < > <i>sessionId</i>	Retorna as informações da sessão.
GET /sessions/ < >/state <i>sessionId</i>	Retorna o estado da sessão.
EXCLUIR /sessions/ < > <i>sessionId</i>	Interrompe e exclui a sessão.
GET /sessions/ < >/declarações <i>sessionId</i>	Retorna todas as instruções em uma sessão.
POST /sessions/ < >/declarações <i>sessionId</i>	Executa de uma instrução em uma sessão.

Operação de API	Description
GET /sessions/ < >/declarações/< > <i>sessionId statementId</i>	Retorna os detalhes da instrução especificada em uma sessão.
POST /sessions/ < >/declarações/< >/cancelar <i>sessionId statementId</i>	Cancela a instrução especificada nesta sessão.

Envio de solicitações ao endpoint do Apache Livy

Você também pode enviar solicitações diretamente ao endpoint do Apache Livy de um cliente HTTP. Isso permite executar remotamente o código para casos de uso fora de um caderno.

Antes de começar a enviar solicitações ao endpoint, certifique-se de que instalou as seguintes bibliotecas:

```
pip3 install botocore awscrt requests
```

Confira este exemplo de script Python para enviar solicitações HTTP diretamente a um endpoint:

```
from botocore import crt
import requests
from botocore.awsrequest import AWSRequest
from botocore.credentials import Credentials
import botocore.session
import json, pprint, textwrap

endpoint = 'https://<application_id>.livy.emr-serverless-
services.<AWS_REGION>.amazonaws.com'
headers = {'Content-Type': 'application/json'}

session = botocore.session.Session()
signer = crt.auth.CrtS3SigV4Auth(session.get_credentials(), 'emr-serverless',
'<AWS_REGION>')

### Create session request

data = {'kind': 'pyspark', 'heartbeatTimeoutInSecond': 60, 'conf': { 'emr-
serverless.session.executionRoleArn': 'arn:aws:iam::123456789012:role/role1'}}
```

```
request = AWSRequest(method='POST', url=endpoint + "/sessions", data=json.dumps(data),
headers=headers)

request.context["payload_signing_enabled"] = False

signer.add_auth(request)

prepped = request.prepare()

r = requests.post(prepped.url, headers=prepped.headers, data=json.dumps(data))

pprint.pprint(r.json())

### List Sessions Request

request = AWSRequest(method='GET', url=endpoint + "/sessions", headers=headers)

request.context["payload_signing_enabled"] = False

signer.add_auth(request)

prepped = request.prepare()

r2 = requests.get(prepped.url, headers=prepped.headers)
pprint.pprint(r2.json())

### Get session state

session_url = endpoint + r.headers['location']

request = AWSRequest(method='GET', url=session_url, headers=headers)

request.context["payload_signing_enabled"] = False

signer.add_auth(request)

prepped = request.prepare()

r3 = requests.get(prepped.url, headers=prepped.headers)

pprint.pprint(r3.json())
```

```
### Submit Statement

data = {
    'code': "1 + 1"
}

statements_url = endpoint + r.headers['location'] + "/statements"

request = AWSRequest(method='POST', url=statements_url, data=json.dumps(data),
    headers=headers)

request.context["payload_signing_enabled"] = False

signer.add_auth(request)

prepped = request.prepare()

r4 = requests.post(prepped.url, headers=prepped.headers, data=json.dumps(data))

pprint.pprint(r4.json())

### Check statements results

specific_statement_url = endpoint + r4.headers['location']

request = AWSRequest(method='GET', url=specific_statement_url, headers=headers)

request.context["payload_signing_enabled"] = False

signer.add_auth(request)

prepped = request.prepare()

r5 = requests.get(prepped.url, headers=prepped.headers)

pprint.pprint(r5.json())

### Delete session

session_url = endpoint + r.headers['location']
```

```
request = AWSRequest(method='DELETE', url=session_url, headers=headers)

request.context["payload_signing_enabled"] = False

signer.add_auth(request)

prepped = request.prepare()

r6 = requests.delete(prepped.url, headers=prepped.headers)

pprint.pprint(r6.json())
```

Kernel do Sparkmagic

Antes de instalar o sparkmagic, verifique se você configurou AWS as credenciais na instância em que deseja instalar o sparkmagic

1. Instale o Sparkmagic seguindo as [etapas de instalação](#). Observe que você só precisa executar as quatro primeiras etapas.
2. O kernel sparkmagic oferece suporte a autenticadores personalizados, para que você possa integrar um autenticador ao kernel sparkmagic para que cada solicitação seja assinada. SIGv4
3. Instale o autenticador personalizado do EMR Sem Servidor.

```
pip install emr-serverless-customauth
```

4. Agora, forneça o caminho para o autenticador personalizado e o URL do endpoint do Apache Livy no arquivo JSON de configuração do Sparkmagic. Use o comando a seguir para abrir o arquivo de configuração.

```
vim ~/.sparkmagic/config.json
```

Veja a seguir um exemplo de arquivo config.json.

```
{
  "kernel_python_credentials" : {
    "username": "",
    "password": "",
    "url": "https://<application-id>.livy.emr-serverless-
services.<AWS_REGION>.amazonaws.com",
    "auth": "Custom_Auth"
```

```

    },

    "kernel_scala_credentials" : {
      "username": "",
      "password": "",
      "url": "https://<application-id>.livy.emr-serverless-
services.<AWS_REGION>.amazonaws.com",
      "auth": "Custom_Auth"
    },
    "authenticators": {
      "None": "sparkmagic.auth.customauth.Authenticator",
      "Basic_Access": "sparkmagic.auth.basic.Basic",
      "Custom_Auth":
"emr_serverless_customauth.customauthenticator.EMRServerlessCustomSigV4Signer"
    },
    "livy_session_startup_timeout_seconds": 600,
    "ignore_ssl_errors": false
  }
}

```

5. Inicie o Jupyter Lab. Ele deve usar a autenticação personalizada que você configurou na última etapa.
6. Em seguida, você pode executar os comandos a seguir do caderno e seu código para começar.

```
%info //Returns the information about the current sessions.
```

```

%%configure -f //Configure information specific to a session. We supply
executionRoleArn in this example. Change it for your use case.
{
  "driverMemory": "4g",
  "conf": {
    "emr-serverless.session.executionRoleArn":
"arn:aws:iam::123456789012:role/JobExecutionRole"
  }
}

```

```
<your code>//Run your code to start the session
```

Internamente, cada instrução chama cada uma das operações da API do Apache Livy por meio do URL configurado do endpoint do Apache Livy. Em seguida, você pode gravar suas instruções de acordo com o caso de uso.

Considerações

Considere as informações a seguir ao executar workloads interativas por meio de endpoints do Apache Livy.

- O EMR Sem Servidor mantém o isolamento em nível de sessão usando a entidade principal do chamador. A entidade principal do chamador que cria a sessão é o único que pode acessá-la. Para um isolamento mais granular, configure uma identidade de origem ao assumir credenciais. Nesse caso, o EMR Sem Servidor impõe o isolamento em nível de sessão com base na entidade principal do chamador e na identidade de origem. Para obter mais informações sobre a identidade de origem, consulte [Monitorar e controlar ações realizadas com perfis assumidos](#).
- Os endpoints do Apache Livy são compatíveis com as versões 6.14.0 e superiores do EMR Sem Servidor.
- Os endpoints do Apache Livy são compatíveis somente com o mecanismo Apache Spark.
- Os endpoints Apache Livy são compatíveis com Scala Spark e PySpark
- Por padrão, `autoStopConfig` está habilitado nas aplicações. Isso significa que as aplicações são encerradas após 15 minutos de inatividade. Você pode alterar essa configuração como parte da solicitação `create-application` ou `update-application`.
- Você pode executar até 25 sessões simultâneas em uma única aplicação habilitada para endpoint do Apache Livy.
- Para obter a melhor experiência de inicialização, sugerimos configurar a capacidade pré-inicializada para drivers e executores.
- Você deve iniciar manualmente a aplicação antes de se conectar ao endpoint do Apache Livy.
- Você deve ter uma cota de serviço de vCPU suficiente Conta da AWS para executar cargas de trabalho interativas com o endpoint Apache Livy. Sugerimos pelo menos 24 vCPUs.
- O tempo limite padrão de uma sessão do Apache Livy é de uma hora. Se você não executar instruções por uma hora, o Apache Livy excluirá a sessão e liberará o driver e os executores. A partir da versão `emr-7.8.0`, esse valor pode ser definido especificando o parâmetro `ttl` como parte da solicitação `/sessions POST` do Livy, por exemplo, 2h (horas), 120m (minutos), 7200s(segundos), 7200000ms (milissegundos).

Note

Essa configuração não pode ser alterada antes do `emr-7.8.0`. A seguir, apresentamos um exemplo de corpo de solicitação `POST /sessions`.

```
{
  "kind": "pyspark",
  "heartbeatTimeoutInSeconds": 60,
  "conf": {
    "emr-serverless.session.executionRoleArn": "executionRoleArn"
  },
  "ttl": "2h"
}
```

- A partir da versão emr-7.8.0 do Amazon EMR para aplicativos com controle de acesso refinado via LakeFormation enabled, a configuração pode ser desativada por sessão. Para obter mais informações sobre como habilitar o controle de acesso refinado para uma aplicação do EMR Sem Servidor, consulte [Métodos para controle de acesso refinado](#).

Note

O Lake Formation não pode ser habilitado para uma Sessão quando não foi habilitado para uma aplicação. A seguir, apresentamos um exemplo de corpo de solicitação POST / sessions.

```
{
  "kind": "pyspark",
  "heartbeatTimeoutInSeconds": 60,
  "conf": {
    "emr-serverless.session.executionRoleArn": "executionRoleArn"
  },
  "spark.emr-serverless.lakeformation.enabled" : "false"
}
```

- Somente sessões ativas podem interagir com um endpoint do Apache Livy. Depois que a sessão for concluída, cancelada ou encerrada, você não poderá acessá-la por meio do endpoint do Apache Livy.

Registro em log e monitoramento

O monitoramento é uma parte importante para manter a confiabilidade, a disponibilidade e a performance de aplicações e trabalhos do Amazon Sem Servidor. Você deve coletar dados de monitoramento de todas as partes de suas soluções EMR Sem Servidor para que falhas multiponto possam ser depuradas mais facilmente, caso ocorram.

Tópicos

- [Armazenamento de logs](#)
- [Logs alternados](#)
- [Criptografia de logs](#)
- [Configuração das propriedades do Apache Log4j2 para o Amazon EMR Sem Servidor](#)
- [Monitoramento do EMR Sem Servidor](#)
- [Automatizando o EMR sem servidor com Amazon EventBridge](#)

Armazenamento de logs

Para monitorar o progresso do seu trabalho no EMR Sem Servidor e solucionar falhas de trabalho, escolha como o EMR Sem Servidor armazena e veicula os logs de aplicações. Ao enviar uma execução de trabalho, especifique armazenamento gerenciado, Amazon S3 e Amazon CloudWatch como suas opções de registro.

Com CloudWatch, especifique os tipos e locais de registro que você deseja usar ou aceite os tipos e locais padrão. Para obter mais informações sobre CloudWatch registros, consulte [the section called “Amazon CloudWatch”](#). Com o armazenamento gerenciado e os logs do S3, a tabela a seguir lista os locais de log e a disponibilidade da interface do usuário que você pode esperar se escolher [armazenamento gerenciado](#), [buckets do Amazon S3](#) ou ambos.

Opção	Logs de eventos	Logs de contêineres	Interfaces de usuário de aplicações
Armazenamento gerenciado	Armazenamento gerenciado	Armazenamento gerenciado	Compatível

Opção	Logs de eventos	Logs de contêineres	Interfaces de usuário de aplicações
Armazenamento gerenciado e bucket do S3	Armazenado em ambos os locais	Armazenado no bucket do S3	Compatível
Bucket do Amazon S3.	Armazenado no bucket do S3	Armazenado no bucket do S3	Sem suporte ¹

¹ Sugerimos manter a opção Armazenamento gerenciado selecionada. Caso contrário, você não poderá usar o aplicativo integrado UIs.

Registro em log do EMR Sem Servidor com armazenamento gerenciado

Por padrão, o EMR Sem Servidor armazena logs de aplicações com segurança no armazenamento gerenciado do Amazon EMR por no máximo 30 dias.

Note

Se você desativar a opção padrão, o Amazon EMR não poderá solucionar problemas de trabalhos em seu nome. Exemplo: você não pode acessar a Spark-UI a partir do console do ERM sem Servidor.

Para desativar essa opção no EMR Studio, desmarque AWS a caixa de seleção Permitir reter registros por 30 dias na seção Configurações adicionais da página Enviar trabalho.

Para desativar essa opção no AWS CLI, use a `managedPersistenceMonitoringConfiguration` configuração ao enviar uma execução de trabalho.

```
{
  "monitoringConfiguration": {
    "managedPersistenceMonitoringConfiguration": {
      "enabled": false
    }
  }
}
```

Se a sua aplicação do EMR Sem Servidor estiver em uma sub-rede privada com endpoints de VPC para o Amazon S3 e você anexar uma política de endpoint para controlar o acesso, adicione as seguintes permissões para que o EMR Sem Servidor armazene e forneça logs de aplicações. Substitua Resource pelos buckets AppInfo da tabela de regiões disponíveis em [Exemplos de políticas para sub-redes privadas que acessam o Amazon S3](#).

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "EMRServerlessManagedLogging",
      "Effect": "Allow",
      "Action": [
        "s3:PutObject",
        "s3:PutObjectAcl"
      ],
      "Resource": [
        "arn:aws:s3:::prod.us-east-1.appinfo.src",
        "arn:aws:s3:::prod.us-east-1.appinfo.src/*"
      ],
      "Condition": {
        "StringEquals": {
          "aws:PrincipalServiceName": "emr-serverless.amazonaws.com",
          "aws:SourceVpc": "vpc-12345678"
        }
      }
    }
  ]
}
```

Além disso, use a chave de condição `aws:SourceVpc` para garantir que a solicitação passe pela VPC à qual o endpoint da VPC está conectado.

Registro em log do EMR Sem Servidor com buckets do Amazon S3

Antes que os trabalhos possam enviar dados de log ao Amazon S3, as permissões apresentadas a seguir devem ser incluídas na política de permissões para o perfil de runtime do trabalho. Substitua *amzn-s3-demo-logging-bucket* pelo nome do bucket de registro em log.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:PutObject"
      ],
      "Resource": [
        "arn:aws:s3:::amzn-s3-demo-bucket/*"
      ],
      "Sid": "AllowS3Putobject"
    }
  ]
}
```

Para configurar um bucket do Amazon S3 para armazenar registros do AWS CLI, use a `s3MonitoringConfiguration` configuração ao iniciar a execução de um trabalho. Para fazer isso, forneça o `--configuration-overrides` a seguir na configuração.

```
{
  "monitoringConfiguration": {
    "s3MonitoringConfiguration": {
      "logUri": "s3://amzn-s3-demo-logging-bucket/logs/"
    }
  }
}
```

Em trabalhos em lotes que não têm novas tentativas habilitadas, o EMR Sem Servidor envia os logs para o seguinte caminho:

```
'/applications/<applicationId>/jobs/<jobId>'
```

Os logs do driver do Spark são armazenados no seguinte caminho pelo EMR Sem Servidor

```
'/applications/<applicationId>/jobs/<jobId>/SPARK_DRIVER/'
```

Os logs do executor do Spark são armazenados no seguinte caminho pelo EMR Sem Servidor

```
'/applications/<applicationId>/jobs/<jobId>/SPARK_EXECUTOR/<EXECUTOR-ID>'
```

O <EXECUTOR-ID> é um número inteiro.

As versões 7.1.0 e posteriores do EMR Sem Servidor oferecem suporte a novas tentativas para trabalhos de streaming e em lote. Se você executar um trabalho com novas tentativas habilitadas, o EMR Sem Servidor adicionará automaticamente um número de tentativas ao prefixo do caminho do log, para que você possa distinguir e rastrear melhor os logs.

```
'/applications/<applicationId>/jobs/<jobId>/attempts/<attemptNumber>/'
```

Registro no EMR Serverless com a Amazon CloudWatch

Ao enviar um trabalho para um aplicativo EMR Serverless, escolha a Amazon CloudWatch como uma opção para armazenar os registros do seu aplicativo. Isso permite que você use recursos de análise de CloudWatch registros, como CloudWatch Logs Insights e Live Tail. Você também pode transmitir registros CloudWatch para outros sistemas, por exemplo, OpenSearch para análise posterior.

O EMR Sem Servidor fornece registro em log em tempo real para logs de drivers. Você pode acessar os registros em tempo real com o recurso de cauda ao CloudWatch vivo ou por meio de comandos de cauda da CloudWatch CLI.

Por padrão, o CloudWatch registro está desativado para o EMR Serverless. Para habilitá-lo, use a configuração em [AWS CLI](#).

Note

A Amazon CloudWatch publica registros em tempo real, portanto, gera mais recursos dos trabalhadores. Se você escolher uma baixa capacidade de trabalhador, o impacto no runtime do trabalho poderá aumentar. Se você ativar o CloudWatch registro, sugerimos que escolha uma capacidade de trabalho maior. Também é possível que a publicação de logs faça controle de utilização se a taxa de transações por segundo (TPS) for muito baixa para PutLogEvents. A configuração de CloudWatch limitação é global para todos os serviços, incluindo o EMR Serverless. Para obter mais informações, consulte [Como determino a limitação em meus registros? CloudWatch](#) em AWS re:post.

Permissões necessárias para fazer login com CloudWatch

Antes que seus trabalhos possam enviar dados de log para a Amazon CloudWatch, inclua as seguintes permissões na política de permissões para a função de tempo de execução do trabalho.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "logs:DescribeLogGroups"
      ],
      "Resource": [
        "arn:aws:logs:*:123456789012:*"
      ],
      "Sid": "AllowLOGSDescribeLogGroups"
    },
    {
      "Effect": "Allow",
      "Action": [
        "logs:PutLogEvents",
        "logs:CreateLogGroup",
        "logs:CreateLogStream",
        "logs:DescribeLogStreams"
      ],
      "Resource": [
        "arn:aws:logs:*:123456789012:log-group:my-log-group-name:*"
      ],
      "Sid": "AllowLOGSPutLogEvents"
    }
  ]
}
```

AWS CLI

Para configurar a Amazon CloudWatch para armazenar registros do EMR Serverless a partir do AWS CLI, use a `cloudWatchLoggingConfiguration` configuração ao iniciar a execução de um trabalho. Para fazer isso, forneça as substituições de configuração a seguir. Opcionalmente, forneça

um nome de grupo de logs, nome de prefixo de fluxo de logs, tipos de log e um ARN de chave de criptografia.

Se você não especificar os valores opcionais, CloudWatch publicará os registros em um grupo de registros padrão/`aws/emr-serverless`, com o fluxo `/applications/applicationId/jobs/jobId/worker-type` de registros padrão.

As versões 7.1.0 e posteriores do EMR Sem Servidor oferecem suporte a novas tentativas para trabalhos de streaming e em lote. Se você habilitou novas tentativas para um trabalho, o EMR Sem Servidor adiciona automaticamente um número de tentativas ao prefixo do caminho do log, para que você possa distinguir e rastrear melhor os logs.

```
'/applications/<applicationId>/jobs/<jobId>/attempts/<attemptNumber>/worker-type'
```

O seguinte demonstra a configuração mínima necessária para ativar o Amazon CloudWatch Logging com as configurações padrão para o EMR Serverless:

```
{
  "monitoringConfiguration": {
    "cloudWatchLoggingConfiguration": {
      "enabled": true
    }
  }
}
```

O exemplo a seguir mostra todas as configurações obrigatórias e opcionais que especificam quando você ativa o Amazon CloudWatch Logging para o EMR Serverless. Os valores de `logTypes` aceitos também estão listados no exemplo a seguir.

```
{
  "monitoringConfiguration": {
    "cloudWatchLoggingConfiguration": {
      "enabled": true, // Required
      "logGroupName": "Example_logGroup", // Optional
      "logStreamNamePrefix": "Example_logStream", // Optional
      "encryptionKeyArn": "key-arn", // Optional
      "logTypes": {
        "SPARK_DRIVER": ["stdout", "stderr"] //List of values
      }
    }
  }
}
```

```
}
```

Por padrão, o EMR Serverless publica somente os registros do driver `stdout` e `stderr` em CloudWatch. Se quiser outros logs, especifique um perfil de contêiner e os tipos de log correspondentes com o campo `logTypes`.

A seguinte lista mostra os tipos de trabalhadores compatíveis que você pode especificar na configuração `logTypes`:

Spark

- `SPARK_DRIVER` : ["STDERR", "STDOUT"]
- `SPARK_EXECUTOR` : ["STDERR", "STDOUT"]

Hive

- `HIVE_DRIVER` : ["STDERR", "STDOUT", "HIVE_LOG", "TEZ_AM"]
- `TEZ_TASK` : ["STDERR", "STDOUT", "SYSTEM_LOGS"]

Logs alternados

O Amazon EMR Sem Servidor pode alternar logs de aplicações e logs de eventos do Spark. A alternância de logs ajuda com o problema de trabalhos de longa execução, gerando grandes arquivos de log que podem ocupar todo o espaço em disco. A alternância de logs ajuda a economizar armazenamento em disco e reduz a quantidade de falhas de trabalho por não haver mais espaço no disco.

A alternância de logs está habilitada por padrão e disponível somente para trabalhos do Spark.

Logs de eventos do Spark

Note

A alternância de logs de eventos do Spark está disponível em todos os rótulos de lançamento do Amazon EMR.

Em vez de gerar um único arquivo de log de eventos, o EMR Sem Servidor alterna o log de eventos em um intervalo de tempo regular e remove os arquivos de log de eventos mais antigos. A alternância de logs não afeta os logs enviados ao bucket do S3.

Logs de aplicação do Spark

Note

A alternância de logs de aplicação do Spark está disponível em todos os rótulos de lançamento do Amazon EMR.

O EMR Sem Servidor também alterna os logs de aplicação do Spark para drivers e executores, como arquivos `stdout` e `stderr`. Você pode acessar os arquivos de log mais recentes escolhendo os links de log no Studio e usando os links da interface de usuário do Live e do servidor de histórico do Spark. Os arquivos de log são as versões truncadas dos logs mais recentes. Para consultar os logs alternados mais antigos, especifique uma localização do Amazon S3 ao armazenar os logs. Consulte [Registro em log para o EMR Sem Servidor com buckets do S3](#) para obter mais informações.

Você encontrará os arquivos de log mais recentes no local a seguir. O EMR Sem Servidor atualiza os arquivos a cada 15 segundos. Esses arquivos podem variar de 0 MB a 128 MB.

```
<example-S3-logUri>/applications/<application-id>/jobs/<job-id>/SPARK_DRIVER/stderr.gz
```

O local a seguir contém os arquivos alternados mais antigos. Cada arquivo tem 128 MB.

```
<example-S3-logUri>/applications/<application-id>/jobs/<job-id>/SPARK_DRIVER/archived/  
stderr_<index>.gz
```

O mesmo comportamento também se aplica aos executores do Spark. Essa alteração é aplicável somente ao registro em log do S3. A rotação de registros não introduz nenhuma alteração nos fluxos de registros enviados para a Amazon CloudWatch.

As versões 7.1.0 e posteriores do EMR Sem Servidor oferecem suporte a novas tentativas para trabalhos de streaming e em lote. Se você habilitou novas tentativas no trabalho, o EMR Sem Servidor adiciona um prefixo ao caminho de log para esses trabalhos, para que você possa rastrear e distinguir melhor os logs uns dos outros. Esse caminho contém todos os logs alternados.

```
'/applications/<applicationId>/jobs/<jobId>/attempts/<attemptNumber>/'.
```

Criptografia de logs

Criptografia de logs do EMR Sem Servidor com armazenamento gerenciado

Para criptografar logs no armazenamento gerenciado com sua própria chave do KMS, use a configuração `managedPersistenceMonitoringConfiguration` ao enviar uma execução de trabalho.

```
{
  "monitoringConfiguration": {
    "managedPersistenceMonitoringConfiguration": {
      "encryptionKeyArn": "key-arn"
    }
  }
}
```

Criptografia de logs do EMR Sem Servidor com buckets do Amazon S3

Para criptografar logs no bucket do Amazon S3 com sua própria chave do KMS, use a configuração `s3MonitoringConfiguration` ao enviar uma execução de trabalho.

```
{
  "monitoringConfiguration": {
    "s3MonitoringConfiguration": {
      "logUri": "s3://amzn-s3-demo-logging-bucket/logs/",
      "encryptionKeyArn": "key-arn"
    }
  }
}
```

Criptografando registros sem servidor do EMR com a Amazon CloudWatch

Para criptografar registros na Amazon CloudWatch com sua própria chave KMS, use a configuração `cloudWatchLoggingConfiguration` ao enviar uma execução de trabalho.

```
{
  "monitoringConfiguration": {
    "cloudWatchLoggingConfiguration": {
```

```

        "enabled": true,
        "encryptionKeyArn": "key-arn"
    }
}
}

```

Permissões necessárias para criptografia de logs

Nesta seção

- [Permissões obrigatórias do usuário](#)
- [Permissões de chave de criptografia para Amazon S3 e armazenamento gerenciado](#)
- [Permissões de chave de criptografia para a Amazon CloudWatch](#)

Permissões obrigatórias do usuário

O usuário que envia o trabalho ou visualiza os registros ou o aplicativo UIs deve ter permissões para usar a chave. Você pode especificar as permissões na política de chave do KMS ou na política do IAM para o usuário, grupo ou perfil. Se o usuário que envia o trabalho não tiver as permissões da chave do KMS, o EMR Sem Servidor rejeitará o envio da execução do trabalho.

Exemplo de política de chaves

A seguinte política de chave fornece permissões para `kms:GenerateDataKey` e `kms:Decrypt`:

```

{
  "Effect": "Allow",
  "Principal": {
    "AWS": "arn:aws:iam::111122223333:user/user-name"
  },
  "Action": [
    "kms:GenerateDataKey",
    "kms:Decrypt"
  ],
  "Resource": "*"
}

```

Exemplo de política do IAM

A seguinte política do IAM fornece as permissões para `kms:GenerateDataKey` e `kms:Decrypt`:

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "kms:GenerateDataKey",
        "kms:Decrypt"
      ],
      "Resource": [
        "arn:aws:kms:*:123456789012:key/12345678-1234-1234-1234-123456789012"
      ],
      "Sid": "AllowKMSGeneratedatakey"
    }
  ]
}
```

Para iniciar a interface de usuário do Spark ou Tez, conceda aos usuários, grupos ou perfis permissões para acessar a API `emr-serverless:GetDashboardForJobRun` da seguinte forma:

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "emr-serverless:GetDashboardForJobRun"
      ],
      "Resource": [
        "*"
      ],
      "Sid": "AllowEMRSERVERLESSGetdashboardforjobrun"
    }
  ]
}
```

Permissões de chave de criptografia para Amazon S3 e armazenamento gerenciado

Ao criptografar logs com sua própria chave de criptografia no armazenamento gerenciado ou nos buckets do S3, configure as permissões da chave do KMS da forma a seguir.

A entidade principal do `emr-serverless.amazonaws.com` deve ter as seguintes permissões na política da chave do KMS:

```
{
  "Effect": "Allow",
  "Principal": {
    "Service": "emr-serverless.amazonaws.com"
  },
  "Action": [
    "kms:Decrypt",
    "kms:GenerateDataKey"
  ],
  "Resource": "*"
  "Condition": {
    "StringLike": {
      "aws:SourceArn": "arn:aws:emr-serverless:region:aws-account-id:/
applications/application-id"
    }
  }
}
```

Como uma prática recomendada de segurança, sugerimos adicionar uma chave de condição `aws:SourceArn` à política de chave do KMS. A chave de condição global do IAM `aws:SourceArn` ajuda a garantir que o EMR Sem Servidor use a chave do KMS somente para o ARN da aplicação.

O perfil de runtime do trabalho deve ter as seguintes permissões na política do IAM:

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "kms:GenerateDataKey",
```

```

    "kms:Decrypt"
  ],
  "Resource": [
    "arn:aws:kms:*:123456789012:key/12345678-1234-1234-1234-123456789012"
  ],
  "Sid": "AllowKMSGeneratedatakey"
}
]
}

```

Permissões de chave de criptografia para a Amazon CloudWatch

Para associar o ARN da chave do KMS ao seu grupo de logs, use a política do IAM a seguir no perfil de runtime do trabalho.

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "logs:AssociateKmsKey"
      ],
      "Resource": [
        "arn:aws:logs:*:123456789012:log-group:my-log-group-name:*"
      ],
      "Sid": "AllowLOGSAssociatekmskey"
    }
  ]
}

```

Configure a política de chaves do KMS para conceder permissões de KMS à Amazon: CloudWatch

JSON

```

{
  "Version": "2012-10-17",

```

```
"Id": "key-default-1",
"Statement": [
  {
    "Effect": "Allow",
    "Action": [
      "kms:Decrypt",
      "kms:GenerateDataKey"
    ],
    "Resource": [
      "*"
    ],
    "Condition": {
      "ArnLike": {
        "kms:EncryptionContext:aws:logs:arn": "arn:aws:logs:*:123456789012:*"
      }
    },
    "Sid": "AllowKMSDecrypt"
  }
]
}
```

Configuração das propriedades do Apache Log4j2 para o Amazon EMR Sem Servidor

Esta página descreve como configurar propriedades personalizadas do [Apache Log4j 2.x](#) para trabalhos do EMR Sem Servidor em StartJobRun. Se quiser configurar as classificações do Log4j no nível da aplicação, consulte [Configuração padrão de aplicações do EMR Sem Servidor](#).

Configuração das propriedades do Log4j2 do Spark para o Amazon EMR Sem Servidor

Com as versões 6.8.0 e posteriores do Amazon EMR, você pode personalizar as propriedades do [Apache Log4j 2.x](#) e especificar configurações de log refinadas. Isso simplifica a solução de problemas dos trabalhos do Spark no EMR Sem Servidor. Para configurar essas propriedades, use as classificações `spark-driver-log4j2` e `spark-executor-log4j2`.

Tópicos

- [Classificações do Log4j2 para o Spark](#)
- [Exemplo de configuração do Log4j2 para o Spark](#)

- [Log4j2 em exemplos de trabalhos do Spark](#)
- [Considerações sobre o Log4j2 para o Spark](#)

Classificações do Log4j2 para o Spark

Para personalizar as configurações de log do Spark, use as classificações a seguir com [applicationConfiguration](#). Para configurar as propriedades do Log4j 2.x, use as [properties](#) a seguir.

spark-driver-log4j2

Essa classificação define os valores no arquivo `log4j2.properties` do driver.

spark-executor-log4j2

Essa classificação define os valores no arquivo `log4j2.properties` do executor.

Exemplo de configuração do Log4j2 para o Spark

O exemplo a seguir mostra como enviar um trabalho do Spark com `applicationConfiguration` para personalizar as configurações do Log4j2 para o driver e o executor do Spark.

Para configurar as classificações do Log4j no nível da aplicação, em vez de quando você envia o trabalho, consulte [Configuração padrão de aplicações do EMR Sem Servidor](#).

```
aws emr-serverless start-job-run \  
  --application-id application-id \  
  --execution-role-arn job-role-arn \  
  --job-driver '{  
    "sparkSubmit": {  
      "entryPoint": "/usr/lib/spark/examples/jars/spark-examples.jar",  
      "entryPointArguments": ["1"],  
      "sparkSubmitParameters": "--class org.apache.spark.examples.SparkPi --conf  
spark.executor.cores=4 --conf spark.executor.memory=20g --conf spark.driver.cores=4 --  
conf spark.driver.memory=8g --conf spark.executor.instances=1"  
    }  
  }'  
  --configuration-overrides '{  
    "applicationConfiguration": [  
      {  
        "classification": "spark-driver-log4j2",
```

```

        "properties": {
            "rootLogger.level": "error", // will only display Spark error logs
            "logger.IdentifierForClass.name": "classpath for setting logger",
            "logger.IdentifierForClass.level": "info"
        }
    },
    {
        "classification": "spark-executor-log4j2",
        "properties": {
            "rootLogger.level": "error", // will only display Spark error logs
            "logger.IdentifierForClass.name": "classpath for setting logger",
            "logger.IdentifierForClass.level": "info"
        }
    }
]
}'

```

Log4j2 em exemplos de trabalhos do Spark

Os exemplos de código a seguir demonstram como criar uma aplicação do Spark enquanto você inicializa uma configuração personalizada do Log4j2 para a aplicação.

Python

Example- Usando o Log4j2 para um trabalho do Spark com Python

```

import os
import sys

from pyspark import SparkConf, SparkContext
from pyspark.sql import SparkSession

app_name = "PySparkApp"
if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .appName(app_name)\
        .getOrCreate()

    spark.sparkContext._conf.getAll()
    sc = spark.sparkContext
    log4jLogger = sc._jvm.org.apache.log4j

```

```
LOGGER = log4jLogger.LogManager.getLogger(app_name)

LOGGER.info("pyspark script logger info")
LOGGER.warn("pyspark script logger warn")
LOGGER.error("pyspark script logger error")

// your code here

spark.stop()
```

Para personalizar o Log4j2 do driver ao executar um trabalho do Spark, use a seguinte configuração:

```
{
  "classification": "spark-driver-log4j2",
  "properties": {
    "rootLogger.level": "error", // only display Spark error logs
    "logger.PySparkApp.level": "info",
    "logger.PySparkApp.name": "PySparkApp"
  }
}
```

Scala

Example- Usando o Log4j2 para um trabalho do Spark com Scala

```
import org.apache.log4j.Logger
import org.apache.spark.sql.SparkSession

object ExampleClass {
  def main(args: Array[String]): Unit = {
    val spark = SparkSession
      .builder
      .appName(this.getClass.getName)
      .getOrCreate()

    val logger = Logger.getLogger(this.getClass);
    logger.info("script logging info logs")
    logger.warn("script logging warn logs")
    logger.error("script logging error logs")

    // your code here
    spark.stop()
  }
}
```

```
}  
}
```

Para personalizar o Log4j2 do driver ao executar um trabalho do Spark, use a seguinte configuração:

```
{  
  "classification": "spark-driver-log4j2",  
  "properties": {  
    "rootLogger.level": "error", // only display Spark error logs  
    "logger.ExampleClass.level": "info",  
    "logger.ExampleClass.name": "ExampleClass"  
  }  
}
```

Considerações sobre o Log4j2 para o Spark

As seguintes propriedades do Log4j2.x não são configuráveis para processos do Spark:

- `rootLogger.appenderRef.stdout.ref`
- `appender.console.type`
- `appender.console.name`
- `appender.console.target`
- `appender.console.layout.type`
- `appender.console.layout.pattern`

[Para obter informações detalhadas sobre as propriedades do Log4j2.x que configuram, consulte o arquivo em `log4j2.properties.template` GitHub](#)

Monitoramento do EMR Sem Servidor

Esta seção aborda as maneiras para monitorar aplicações e trabalhos do Amazon EMR Sem Servidor.

Tópicos

- [Monitoramento de aplicações e trabalhos do EMR Sem Servidor](#)
- [Monitoramento de métricas do Spark com o Amazon Managed Service para Prometheus](#)

- [Métricas de uso do EMR Sem Servidor](#)

Monitoramento de aplicações e trabalhos do EMR Sem Servidor

Com CloudWatch as métricas da Amazon para EMR Serverless, você pode receber CloudWatch métricas de 1 minuto e acessar CloudWatch painéis para acessar near-real-time as operações e o desempenho de seus aplicativos EMR Serverless.

O EMR Serverless envia métricas a cada minuto. CloudWatch O EMR Serverless emite essas métricas no nível do aplicativo, bem como no cargo, no tipo de funcionário e nos níveis. `capacity-allocation-type`

Para começar, use o modelo de CloudWatch painel do EMR Serverless fornecido no repositório do [EMR](#) Serverless e implante-o. [GitHub](#)

Note

As [workloads interativas do EMR Sem Servidor](#) têm apenas o monitoramento em nível de aplicação habilitado e uma nova dimensão de tipo de trabalhador, `Spark_Kernel`. Para monitorar e depurar suas workloads interativas, acesse os logs e a interface do usuário do Apache Spark no [Workspace do EMR Studio](#).

Monitoramento de métricas

Important

Estamos reestruturando nossa exibição de métricas para adicionar `ApplicationName` e `JobName` como dimensões. Para a versão 7.10 e posteriores, as métricas mais antigas não serão mais atualizadas. Para versões do EMR abaixo de 7.10, as métricas mais antigas ainda estão disponíveis.

Dimensões atuais

A tabela abaixo descreve as dimensões do EMR Sem Servidor disponíveis no namespace `AWS/EMR Serverless`.

Dimensões para métricas do EMR Sem Servidor

Dimensão	Description	
ApplicationId	Filtros para todas as métricas de uma aplicação EMR Sem Servidor usando o ID da aplicação.	
ApplicationName	Filtros para todas as métricas de uma aplicação EMR Sem Servidor usando o nome. Se o nome não for fornecido ou contiver caracteres não ASCII, ele será publicado como [Unspecified].	
JobId	Filtros para todas as métricas da execução de um ID de execução de trabalho do EMR Sem Servidor.	
JobName	Filtros para todas as métricas de um trabalho EMR Sem Servidor executado usando o nome. Se o nome não for fornecido ou contiver caracteres não ASCII, ele será publicado como [Unspecified].	
WorkerType	Filtros para todas as métricas de um determinado tipo de trabalhador. Por exemplo, você pode filtrar por SPARK_DRIVER e SPARK_EXECUTORS para trabalhos do Spark.	

Dimensão	Description
CapacityAllocationType	Filtros para todas as métricas de um determinado tipo de alocação de capacidade. Por exemplo, você pode filtrar por <code>PreInitCapacity</code> para capacidade pré-inicializada e <code>OnDemandCapacity</code> para todo o resto.

Monitoramento em nível de aplicações

Você pode monitorar o uso da capacidade no nível do aplicativo EMR Serverless com as métricas da Amazon CloudWatch. Você também pode configurar um único monitor para monitorar o uso da capacidade do aplicativo em um CloudWatch painel.

Métricas de aplicações do EMR Sem Servidor

Métrica	Description	Unidade	Dimensão
MaxCPUAllowed	O máximo de CPU permitido para a aplicação.	vCPU	ApplicationId , ApplicationName
MaxMemoryAllowed	A memória máxima em GB permitida para a aplicação.	Gigabytes (GB)	ApplicationId , ApplicationName
MaxStorageAllowed	O armazenamento máximo em GB permitido para a aplicação.	Gigabytes (GB)	ApplicationId , ApplicationName
CPUAllocated	Os números totais de v CPUs alocados.	vCPU	ApplicationId , ApplicationName , WorkerType

Métrica	Description	Unidade	Dimensão
			ApplicationId , CapacityAllocationType
IdleWorkerCount	O número total de trabalhadores ociosos.	Contagem	ApplicationId , ApplicationName , WorkerType , CapacityAllocationType
MemoryAllocated	A memória total em GB alocada.	Gigabytes (GB)	ApplicationId , ApplicationName , WorkerType , CapacityAllocationType
PendingCreationWorkerCount	O número total de trabalhadores pendentes de criação.	Contagem	ApplicationId , ApplicationName , WorkerType , CapacityAllocationType
RunningWorkerCount	O número total de trabalhadores em uso pela aplicação.	Contagem	ApplicationId , ApplicationName , WorkerType , CapacityAllocationType
StorageAllocated	O armazenamento total em disco em GB alocado.	Gigabytes (GB)	ApplicationId , ApplicationName , WorkerType , CapacityAllocationType

Métrica	Description	Unidade	Dimensão
TotalWorkerCount	O número total de trabalhadores disponíveis.	Contagem	ApplicationId , ApplicationName , WorkerType , CapacityAllocationType

Monitoramento no nível do trabalho

O Amazon EMR Sem Servidor envia as métricas de nível de trabalho a seguir a cada minuto ao Amazon CloudWatch . Você pode acessar os valores das métricas para execuções de trabalhos agregadas por estado de execução de trabalhos. A unidade para cada uma das métricas é contagem.

Métricas de nível de trabalho do EMR Sem Servidor

Métrica	Description	Dimensão
SubmittedJobs	O número de trabalhos no estado Enviado.	ApplicationId , ApplicationName
PendingJobs	O número de trabalhos em um estado Pendente.	ApplicationId , ApplicationName
ScheduledJobs	O número de trabalhos em um estado Programado.	ApplicationId , ApplicationName
RunningJobs	O número de trabalhos em um estado Em execução.	ApplicationId , ApplicationName
SuccessJobs	O número de trabalhos em um estado Com êxito.	ApplicationId , ApplicationName
FailedJobs	O número de trabalhos em um estado de Falha.	ApplicationId , ApplicationName
CancellingJobs	O número de trabalhos em um estado de Cancelamento.	ApplicationId , ApplicationName

Métrica	Description	Dimensão
CancelledJobs	O número de trabalhos em um estado Cancelado.	ApplicationId , ApplicationName

Você pode monitorar métricas específicas do mecanismo para executar e concluir trabalhos do EMR Serverless com um aplicativo específico do mecanismo. Uls Quando você acessa a interface do usuário para um trabalho em execução, a interface do usuário da aplicação ao vivo é exibida com atualizações em tempo real. Quando você acessa a interface do usuário para um trabalho concluído, a interface do usuário persistente da aplicação é exibida.

Execução de trabalhos

Para trabalhos do EMR Sem Servidor em execução, acesse uma interface em tempo real que fornece métricas específicas do mecanismo. Você pode usar a interface do usuário do Apache Spark ou a interface do usuário do Hive Tez para monitorar e depurar trabalhos. Para acessá-los UIs, use o console do EMR Studio ou solicite um endpoint de URL seguro com o AWS Command Line Interface

Trabalhos concluídos

Para trabalhos concluídos do EMR Sem Servidor, use o Spark History Server ou a interface de usuário persistente do Hive Tez para acessar detalhes, estágios, tarefas e métricas das execuções de trabalhos do Spark ou do Hive. Para acessá-los UIs, use o console do EMR Studio ou solicite um endpoint de URL seguro com o AWS Command Line Interface

Monitoramento em nível de trabalhador para trabalhos

O Amazon EMR Serverless envia as seguintes métricas de nível de funcionário que estão disponíveis no `AWS/EMRServerless` namespace e no grupo de métricas para a Amazon. `Job Worker Metrics CloudWatch` O EMR Serverless coleta pontos de dados de trabalhadores individuais durante a execução do trabalho no nível do cargo, no tipo de trabalhador e no nível. `capacity-allocation-type` Você pode usar `ApplicationId` como uma dimensão para monitorar vários trabalhos que pertencem à mesma aplicação.

Note

Para visualizar o total de CPU e memória usadas por um trabalho sem servidor do EMR ao visualizar as métricas no CloudWatch console da Amazon, use a estatística como soma e o período como 1 minuto.

Métricas do EMR Sem Servidor em nível de trabalhador

Métrica	Description	Unidade	Dimensão
WorkerCpuAllocated	O número total de núcleos de vCPU alocados para trabalhadores em uma execução de trabalho.	vCPU	JobId, JobName, ApplicationId, ApplicationName, WorkerType, e CapacityAllocationType
WorkerCpuUsed	O número total de núcleos de vCPU utilizados pelos trabalhadores em uma execução de trabalho.	vCPU	JobId, JobName, ApplicationId, ApplicationName, WorkerType, e CapacityAllocationType
WorkerMemoryAllocated	A memória total em GB alocada para trabalhadores em uma	Gigabytes (GB)	JobId, JobName, ApplicationId, ApplicationName, WorkerType, e CapacityAllocationType

Métrica	Description	Unidade	Dimensão
	execução de trabalho.		onName , WorkerType , e CapacityAllocation Type
WorkerMemoryUsed	A memória total em GB utilizada pelos trabalhadores em uma execução de trabalho.	Gigabytes (GB)	JobId, JobName, ApplicationId , ApplicationName , WorkerType , e CapacityAllocation Type
WorkerEphemeralStorageAllocated	O número de bytes de armazenamento temporário alocados para trabalhadores em uma execução de trabalho.	Gigabytes (GB)	JobId, JobName, ApplicationId , ApplicationName , WorkerType , e CapacityAllocation Type

Métrica	Description	Unidade	Dimensão
<code>WorkerEphemeralStorageUsed</code>	O número de bytes de armazenamento temporário usados pelos trabalhadores em uma execução de trabalho.	Gigabytes (GB)	JobId, JobName, ApplicationId, ApplicationName, WorkerType, e CapacityAllocationType
<code>WorkerStorageReadBytes</code>	O número de bytes lidos do armazenamento por trabalhadores durante uma execução de trabalho.	Bytes	JobId, JobName, ApplicationId, ApplicationName, WorkerType, e CapacityAllocationType
<code>WorkerStorageWriteBytes</code>	O número de bytes gravados no armazenamento por trabalhadores durante uma execução de trabalho.	Bytes	JobId, JobName, ApplicationId, ApplicationName, WorkerType, e CapacityAllocationType

As etapas abaixo descrevem como acessar os vários tipos de métricas.

Console

Para acessar a interface do usuário da aplicação com o console

1. Navegue até a aplicação do EMR Sem Servidor no EMR Studio com as instruções em [Getting started from the console](#).
2. Para acessar aplicativos UIs e registros específicos do mecanismo para um trabalho em execução:
 - a. Escolha um trabalho com um status RUNNING.
 - b. Selecione o trabalho na página de Detalhes da aplicação ou navegue até a página Detalhes do trabalho do seu trabalho.
 - c. No menu suspenso Exibir interface do usuário, escolha Interface do usuário do Spark ou Interface do usuário do Hive Tez para navegar até a interface da aplicação do seu tipo de trabalho.
 - d. Para acessar os logs do mecanismo do Spark, navegue até a guia Executores na interface do usuário do Spark e escolha o link Logs do driver. Para acessar os logs do mecanismo do Hive, escolha o link Logs do DAG apropriado na interface do usuário do Hive Tez.
3. Para acessar aplicativos UIs e registros específicos do motor para um trabalho concluído:
 - a. Escolha um trabalho com um status SUCCESS.
 - b. Selecione o trabalho na página Detalhes da aplicação ou navegue até a página Detalhes do trabalho.
 - c. No menu suspenso Exibir interface do usuário, escolha Servidor de histórico do Spark ou Interface de usuário persistente do Hive Tez para navegar até a interface da aplicação do seu tipo de trabalho.
 - d. Para acessar os logs do mecanismo do Spark, navegue até a guia Executores na interface do usuário do Spark e escolha o link Logs do driver. Para acessar os logs do mecanismo do Hive, escolha o link Logs do DAG apropriado na interface do usuário do Hive Tez.

AWS CLI

Para acessar a interface do usuário do seu aplicativo com o AWS CLI

- Para gerar um URL usado para acessar a interface do usuário da aplicação para trabalhos em execução e concluídos, chame a API `GetDashboardForJobRun`.

```
aws emr-serverless get-dashboard-for-job-run /  
--application-id <application-id> /  
--job-run-id <job-id>
```

O URL gerado é válido por uma hora.

Monitoramento de métricas do Spark com o Amazon Managed Service para Prometheus

Com as versões 7.1.0 e posteriores do Amazon EMR, você pode integrar o EMR Sem Servidor ao Amazon Managed Service for Prometheus para coletar métricas do Apache Spark para trabalhos e aplicações do EMR Sem Servidor. Essa integração está disponível quando você envia um trabalho ou cria um aplicativo usando o AWS console, a API do EMR Serverless ou o AWS CLI

Pré-requisitos

Antes de fornecer as métricas do Spark para o Amazon Managed Service for Prometheus, conclua os pré-requisitos a seguir.

- [Crie um Workspace do Amazon Managed Service para Prometheus](#). Este Workspace serve como um endpoint de ingestão. Anote o URL exibido em Endpoint: URL de gravação remota. Você precisará especificar o URL ao criar sua aplicação do EMR Sem Servidor.
- Para conceder acesso aos seus trabalhos ao Amazon Managed Service for Prometheus para fins de monitoramento, adicione a política a seguir ao perfil de execução de tarefas.

```
{  
  "Sid": "AccessToPrometheus",  
  "Effect": "Allow",  
  "Action": ["aps:RemoteWrite"],  
  "Resource": "arn:aws:aps:<AWS_REGION>:<AWS_ACCOUNT_ID>:workspace/<WORKSPACE_ID>"  
}
```

Configuração

Para usar o AWS console para criar um aplicativo integrado ao Amazon Managed Service for Prometheus

1. Consulte [Getting started with Amazon EMR Sem Servidor](#) para criar uma aplicação.
2. Ao criar uma aplicação, escolha Usar configurações personalizadas e, em seguida, configure-a especificando as informações nos campos que deseja configurar.
3. Em Logs e métricas da aplicação, escolha Fornecer métricas do mecanismo ao Amazon Managed Service for Prometheus e, em seguida, especifique o URL de gravação remota.
4. Especifique todas as outras configurações desejadas e escolha Criar e iniciar aplicação.

Use a API AWS CLI sem servidor do EMR

Você também pode usar a API AWS CLI ou EMR Serverless para integrar seu aplicativo EMR Serverless ao Amazon Managed Service for Prometheus ao executar o ou os comandos. `create-application` `start-job-run`

`create-application`

```
aws emr-serverless create-application \  
--release-label emr-7.1.0 \  
--type "SPARK" \  
--monitoring-configuration '{  
    "prometheusMonitoringConfiguration": {  
        "remoteWriteUrl": "https://aps-workspaces.<AWS_REGION>.amazonaws.com/  
workspaces/<WORKSPACE_ID>/api/v1/remote_write"  
    }  
'
```

`start-job-run`

```
aws emr-serverless start-job-run \  
--application-id <APPLICATION_ID> \  
--execution-role-arn <JOB_EXECUTION_ROLE> \  
--job-driver '{  
    "sparkSubmit": {  
        "entryPoint": "local:///usr/lib/spark/examples/src/main/python/pi.py",  
        "entryPointArguments": ["10000"],
```

```

        "sparkSubmitParameters": "--conf spark.dynamicAllocation.maxExecutors=10"
    }
}' \
--configuration-overrides '{
    "monitoringConfiguration": {
        "prometheusMonitoringConfiguration": {
            "remoteWriteUrl": "https://aps-workspaces.<AWS_REGION>.amazonaws.com/
workspaces/<WORKSPACE_ID>/api/v1/remote_write"
        }
    }
}'

```

Incluir `prometheusMonitoringConfiguration` no comando indica que o EMR Sem Servidor deve executar o trabalho do Spark com um agente que coleta as métricas do Spark e as grava no endpoint `remoteWriteUrl` para o Amazon Managed Service for Prometheus. Em seguida, você pode usar as métricas do Spark no Amazon Managed Service for Prometheus para exibição, alertas e análises.

Propriedades de configuração avançada

O EMR Sem Servidor usa um componente do Spark chamado `PrometheusServlet` para coletar métricas do Spark e traduzir dados de performance em dados compatíveis com o Amazon Managed Service for Prometheus. Por padrão, o EMR Sem Servidor define valores padrão no Spark e analisa as métricas do driver e do executor quando você envia um trabalho usando `PrometheusMonitoringConfiguration`.

A tabela a seguir descreve todas as propriedades que você pode configurar ao enviar um trabalho do Spark que envia métricas ao Amazon Managed Service for Prometheus.

Propriedade do Spark	Valor padrão	Description
<code>spark.metrics.conf</code> <code>.*.sink.prometheus</code> <code>Servlet.class</code>	<code>org.apache.spark.metrics.sink.</code> <code>PrometheusServlet</code>	A classe que o Spark usa para enviar métricas ao Amazon Managed Service for Prometheus. Para substituir o comportamento padrão, especifique sua própria classe personalizada.

Propriedade do Spark	Valor padrão	Description
<code>spark.metrics.conf</code> <code>.*.source.jvm.class</code>	<code>org.apache.spark.metrics.source.JvmSource</code>	A classe que o Spark usa para coletar e enviar métricas cruciais da máquina virtual Java subjacente. Para parar de coletar métricas da JVM, desabilite essa propriedade definindo-a como uma string vazia, como <code>""</code> . Para substituir o comportamento padrão, especifique sua própria classe personalizada.
<code>spark.metrics.conf</code> <code>.driver.sink.prometheusServlet.path</code>	<code>/metrics/prometheus</code>	O URL distinto que o Amazon Managed Service for Prometheus usa para coletar métricas do driver. Para substituir o comportamento padrão, especifique seu próprio caminho. Para parar de coletar métricas do driver, desabilite essa propriedade definindo-a como uma string vazia, como <code>""</code> .
<code>spark.metrics.conf</code> <code>.executor.sink.prometheusServlet.path</code>	<code>/metrics/executor/prometheus</code>	O URL distinto que o Amazon Managed Service for Prometheus usa para coletar métricas do executor. Para substituir o comportamento padrão, especifique seu próprio caminho. Para parar de coletar métricas do executor, desabilite essa propriedade definindo-a como uma string vazia, como <code>""</code> .

Para obter mais informações sobre as métricas do Spark, consulte [Métricas do Apache Spark](#).

Considerações e limitações

Ao usar o Amazon Managed Service for Prometheus para coletar métricas do EMR Sem Servidor, considere as informações e limitações a seguir.

- O suporte para o uso do Amazon Managed Service for Prometheus com o EMR Sem Servidor está disponível somente nas [Regiões da AWS onde o Amazon Managed Service for Prometheus encontra-se disponível ao público geral](#).
- Executar o agente para coletar métricas do Spark no Amazon Managed Service for Prometheus exige mais recursos dos trabalhadores. Se você escolher um trabalhador menor, como um trabalhador de uma vCPU, o runtime do trabalho poderá aumentar.
- O suporte para o uso do Amazon Managed Service for Prometheus com o EMR Sem Servidor está disponível somente nas versões 7.1.0 e posteriores do Amazon EMR.
- O Amazon Managed Service for Prometheus deve ser implantado na mesma conta em que você executa o EMR Sem Servidor para coletar métricas.

Métricas de uso do EMR Sem Servidor

Você pode usar as métricas CloudWatch de uso da Amazon para dar visibilidade aos recursos que sua conta usa. Use essas métricas para visualizar seu uso do serviço em CloudWatch gráficos e painéis.

As métricas de uso do EMR Sem Servidor correspondem ao Service Quotas. Também é possível configurar alarmes que alertem você quando o uso se aproximar de uma cota de serviço. Para obter mais informações, consulte [Service Quotas e CloudWatch alarmes da Amazon](#) no Service Quotas User Guide.

Para obter mais informações sobre as cotas de serviço do EMR Sem Servidor, consulte [Endpoints e cotas do EMR Serverless](#).

Métricas de uso da cota de serviço do EMR Sem Servidor

O EMR Sem servidor publica as métricas de uso da cota de serviço a seguir no namespace AWS/Usage.

Métrica	Description
ResourceCount	O número total dos recursos especificados em execução na conta. O recurso é definido pelas dimensões associadas à métrica.

Dimensões das métricas de uso da cota de serviço do EMR Sem Servidor

Você pode usar as dimensões a seguir para refinar as métricas de uso que o EMR Sem Servidor publica.

Dimensão	Valor	Description
Service	EMR Sem Servidor	O nome do AWS service (Serviço da AWS) que contém o recurso.
Type	Recurso	O tipo de entidade que o EMR Sem Servidor está relatando.
Resource	vCPU	O tipo de recurso que o EMR Sem Servidor está rastreando.
Class	Nenhum	A classe de recurso que o EMR Sem Servidor está rastreando.

Automatizando o EMR sem servidor com Amazon EventBridge

Você pode usar Amazon EventBridge para automatizar Serviços da AWS e responder automaticamente aos eventos do sistema, como problemas de disponibilidade de aplicativos ou alterações de recursos. EventBridge fornece um fluxo quase em tempo real de eventos do sistema que descrevem as mudanças em seus AWS recursos. Você pode escrever regras simples para indicar quais eventos são do seu interesse, e as ações automatizadas a serem tomadas quando um evento corresponder à regra. Com EventBridge, você pode automaticamente:

- Invocar uma função AWS Lambda
- Retransmissão de um evento para o Amazon Kinesis Data Streams
- Ativar uma máquina de AWS Step Functions estado
- Notificação de um tópico do Amazon SNS ou de uma fila do Amazon SQS

Por exemplo, ao usar EventBridge com o EMR Serverless, você pode ativar uma AWS Lambda função quando uma tarefa de ETL for bem-sucedida ou notificar um tópico do Amazon SNS quando uma tarefa de ETL falhar.

O EMR Sem Servidor emite quatro tipos de eventos:

- Eventos de mudança de estado da aplicação: eventos que emitem todas as alterações de estado de uma aplicação. Para obter mais informações sobre estados da aplicação, consulte [Estados da aplicação](#).
- Eventos de mudança de estado de execução de um trabalho: eventos que emitem cada mudança de estado da execução de um trabalho. Para obter mais informações, consulte [Estados de execução de trabalho](#).
- Eventos de repetição da execução de um trabalho: eventos que emitem cada nova tentativa de um trabalho executado nas versões 7.1.0 e posteriores do Amazon EMR Sem Servidor.
- Eventos de atualização de utilização de recursos de trabalho: eventos que emitem atualizações de utilização de recursos para um trabalho executado em intervalos de aproximadamente 30 minutos.

Exemplos de eventos do EMR Serverless EventBridge

Os eventos relatados pelo EMR Sem Servidor têm um valor de `aws.emr-serverless` atribuído a `source`, como nos exemplos a seguir.

Evento de alteração do estado da aplicação

O exemplo de evento a seguir mostra uma aplicação no estado CREATING.

```
{
  "version": "0",
  "id": "9fd3cf79-1ff1-b633-4dd9-34508dc1e660",
  "detail-type": "EMR Serverless Application State Change",
  "source": "aws.emr-serverless",
  "account": "123456789012",
```

```

"time": "2022-05-31T21:16:31Z",
"region": "us-east-1",
"resources": [],
"detail": {
  "applicationId": "00f1cb5c6anuij25",
  "applicationName": "3965ad00-8fba-4932-a6c8-ded32786fd42",
  "arn": "arn:aws:emr-serverless:us-east-1:111122223333:/
applications/00f1cb5c6anuij25",
  "releaseLabel": "emr-6.6.0",
  "state": "CREATING",
  "type": "HIVE",
  "createdAt": "2022-05-31T21:16:31.547953Z",
  "updatedAt": "2022-05-31T21:16:31.547970Z",
  "autoStopConfig": {
    "enabled": true,
    "idleTimeout": 15
  },
  "autoStartConfig": {
    "enabled": true
  }
}
}

```

Eventos de alteração de estado em execuções de trabalhos

O exemplo de evento a seguir mostra uma execução de trabalho que se move de um estado SCHEDULED para RUNNING.

```

{
  "version": "0",
  "id": "00df3ec6-5da1-36e6-ab71-20f0de68f8a0",
  "detail-type": "EMR Serverless Job Run State Change",
  "source": "aws.emr-serverless",
  "account": "123456789012",
  "time": "2022-05-31T21:07:42Z",
  "region": "us-east-1",
  "resources": [],
  "detail": {
    "jobRunId": "00f1cbn5g4bb0c01",
    "applicationId": "00f1982r1uukb925",
    "arn": "arn:aws:emr-serverless:us-east-1:123456789012:/
applications/00f1982r1uukb925/jobruns/00f1cbn5g4bb0c01",
    "releaseLabel": "emr-6.6.0",

```

```

    "state": "RUNNING",
    "previousState": "SCHEDULED",
    "createdBy": "arn:aws:sts::123456789012:assumed-role/
TestRole-402dcef3ad14993c15d28263f64381e4cda34775/6622b6233b6d42f59c25dd2637346242",
    "updatedAt": "2022-05-31T21:07:42.299487Z",
    "createdAt": "2022-05-31T21:07:25.325900Z"
  }
}

```

Evento de repetição de execução de trabalho

Confira a seguir um exemplo de um evento de nova tentativa de execução de trabalho.

```

{
  "version": "0",
  "id": "00df3ec6-5da1-36e6-ab71-20f0de68f8a0",
  "detail-type": "EMR Serverless Job Run Retry",
  "source": "aws.emr-serverless",
  "account": "123456789012",
  "time": "2022-05-31T21:07:42Z",
  "region": "us-east-1",
  "resources": [],
  "detail": {
    "jobRunId": "00f1cbn5g4bb0c01",
    "applicationId": "00f1982r1uukb925",
    "arn": "arn:aws:emr-serverless:us-east-1:123456789012:/
applications/00f1982r1uukb925/jobruns/00f1cbn5g4bb0c01",
    "releaseLabel": "emr-6.6.0",
    "createdBy": "arn:aws:sts::123456789012:assumed-role/
TestRole-402dcef3ad14993c15d28263f64381e4cda34775/6622b6233b6d42f59c25dd2637346242",
    "updatedAt": "2022-05-31T21:07:42.299487Z",
    "createdAt": "2022-05-31T21:07:25.325900Z",
    //Attempt Details
    "previousAttempt": 1,
    "previousAttemptState": "FAILED",
    "previousAttemptCreatedAt": "2022-05-31T21:07:25.325900Z",
    "previousAttemptEndedAt": "2022-05-31T21:07:30.325900Z",
    "newAttempt": 2,
    "newAttemptCreatedAt": "2022-05-31T21:07:30.325900Z"
  }
}

```

Atualização da utilização de recursos de trabalho

O exemplo de evento a seguir mostra a atualização final da utilização de recursos para um trabalho que foi transferido para um estado terminal após a execução.

```
{
  "version": "0",
  "id": "00df3ec6-5da1-36e6-ab71-20f0de68f8a0",
  "detail-type": "EMR Serverless Job Resource Utilization Update",
  "source": "aws.emr-serverless",
  "account": "123456789012",
  "time": "2022-05-31T21:07:42Z",
  "region": "us-east-1",
  "resources": [
    "arn:aws:emr-serverless:us-east-1:123456789012:/applications/00f1982r1uukb925/jobruns/00f1cbn5g4bb0c01"
  ],
  "detail": {
    "applicationId": "00f1982r1uukb925",
    "jobRunId": "00f1cbn5g4bb0c01",
    "attempt": 1,
    "mode": "BATCH",
    "createdAt": "2022-05-31T21:07:25.325900Z",
    "startedAt": "2022-05-31T21:07:26.123Z",
    "calculatedFrom": "2022-05-31T21:07:42.299487Z",
    "calculatedTo": "2022-05-31T21:07:30.325900Z",
    "resourceUtilizationFinal": true,
    "resourceUtilizationForInterval": {
      "vCPUHour": 0.023,
      "memoryGBHour": 0.114,
      "storageGBHour": 0.228
    },
    "billedResourceUtilizationForInterval": {
      "vCPUHour": 0.067,
      "memoryGBHour": 0.333,
      "storageGBHour": 0
    },
    "totalResourceUtilization": {
      "vCPUHour": 0.023,
      "memoryGBHour": 0.114,
      "storageGBHour": 0.228
    },
    "totalBilledResourceUtilization": {
      "vCPUHour": 0.067,
      "memoryGBHour": 0.333,
```

```
    "storageGBHour": 0
  }
}
```

O campo `startedAt` só estará presente no caso de o trabalho ter sido movido para um estado de execução.

Marcando atributos

Atribua seus próprios metadados a cada recurso usando tags para ajudar a gerenciar os recursos do EMR Sem Servidor. Esta seção fornece uma visão geral das funções das tags e mostra como você pode criá-las.

Tópicos

- [O que é uma tag?](#)
- [Marcar recursos](#)
- [Limitações de tags](#)
- [Trabalhando com tags usando a AWS CLI e a API Serverless do Amazon EMR](#)

O que é uma tag?

Uma tag é um rótulo que você atribui a um AWS recurso. Cada tag consiste em uma chave e um valor, ambos definidos por você. As tags permitem que você categorize seus AWS recursos por atributos como propósito, proprietário e ambiente. Caso possua muitos recursos do mesmo tipo, identifique rapidamente um recurso específico com base nas tags atribuídas a ele. Por exemplo, defina um conjunto de tags para as aplicações do Amazon EMR Sem Servidor a fim de ajudar a rastrear o proprietário e o nível de pilha de cada aplicação. Sugerimos planejar um conjunto consistente de chaves de tags para cada tipo de recurso.

Tags não são automaticamente atribuídas aos recursos. Depois de adicionar uma tag a um recurso, modifique o valor da tag ou remova-a do recurso a qualquer momento. As tags não têm significado semântico no Amazon EMR Sem Servidor e são interpretadas estritamente como uma string de caracteres. Se você adicionar uma etiqueta que tenha a mesma chave de uma etiqueta existente nesse recurso, o novo valor substituirá o antigo.

Se você usa o IAM, você pode controlar quais usuários em sua AWS conta têm permissão para gerenciar tags. Para obter exemplos de políticas de controle de acesso por etiquetas, consulte [Políticas para controle de acesso baseado em etiquetas](#).

Marcar recursos

Você pode marcar aplicações novas ou existentes e execuções de trabalhos. Se você estiver usando a API Serverless do Amazon EMR, a ou um AWS SDK AWS CLI, poderá aplicar tags a novos

recursos usando o `tags` parâmetro na ação relevante da API. Você também pode aplicar etiquetas aos recursos usando a ação da API `TagResource`.

Você pode usar algumas ações de criação de recursos para especificar etiquetas para um recurso quando ele for criado. Nesse caso, se as etiquetas não puderem ser aplicadas enquanto o recurso estiver sendo criado, ele não será criado. Esse mecanismo garante que os recursos que você pretende etiquetar na criação sejam criados com as etiquetas especificadas ou não sejam criados. Se você marcar os recursos no momento da criação, não precisará executar os scripts de marcação personalizados após a criação do recurso.

A tabela a seguir descreve os recursos do Amazon EMR Sem Servidor que podem ser marcados.

Recursos do que podem ser marcados

Recurso	Compatível com tags	Compatível com a propagação de tags	Suporta marcação na criação (Amazon EMR Serverless API AWS CLI e SDK) AWS	API para criação (etiquetas podem ser adicionadas durante a criação)
Aplicação	Sim	Não. As tags associadas a uma aplicação não se propagam para execuções de trabalhos enviadas a essa aplicação.	Sim	<code>CreateApplication</code>
Execução de trabalho	Sim	Não	Sim	<code>StartJobRun</code>

Limitações de tags

As seguintes limitações básicas se aplicam às tags:

- Cada recurso pode ter no máximo 50 tags criadas pelo usuário.
- Em todos os recursos, cada chave de tag deve ser exclusiva e possuir apenas um valor.
- O comprimento máximo da chave da é de 128 caracteres Unicode em UTF-8.
- O comprimento máximo do valor da é de 256 caracteres Unicode em UTF-8.
- Os caracteres permitidos são letras, números, espaços representáveis em UTF-8 e os seguintes caracteres: `_ . : / = + - @`.
- Uma chave de etiqueta não pode ser uma string vazia. Um valor de tag pode ser uma string vazia, mas não nula.
- As chaves e valores das tags diferenciam maiúsculas de minúsculas.
- Não use `AWS :` ou qualquer combinação de letras maiúsculas e minúsculas como prefixo para chaves ou valores. Esses são reservados para uso pela AWS .

Trabalhando com tags usando a AWS CLI e a API Serverless do Amazon EMR

Use os AWS CLI comandos a seguir ou as operações de API sem servidor do Amazon EMR para adicionar, atualizar, listar e excluir as tags dos seus recursos.

Comandos da CLI e operações de API para tags

Recurso	Compatível com tags	Compatível com a propagação de tags
Adicione ou substitua uma ou mais tags	<code>tag-resource</code>	<code>TagResource</code>
Lista de tags para um recurso	<code>list-tags-for-resource</code>	<code>ListTagsForResource</code>
Exclua uma ou mais tags	<code>untag-resource</code>	<code>UntagResource</code>

Os exemplos a seguir mostram como marcar ou desmarcar recursos usando AWS CLI.

Tags para aplicações existentes

O comando a seguir marca uma aplicação existente.

```
aws emr-serverless tag-resource --resource-arn resource_ARN --tags team=devs
```

Desmarcação de uma aplicação existente

O comando a seguir exclui uma tag de uma aplicação existente.

```
aws emr-serverless untag-resource --resource-arn resource_ARN --tag-keys tag_key
```

Lista de tags para um recurso

O comando a seguir lista as tags associadas a um recurso existente.

```
aws emr-serverless list-tags-for-resource --resource-arn resource_ARN
```

Tutoriais do EMR Sem Servidor

Esta seção descreve casos de uso comuns quando você trabalha com aplicações do EMR Sem Servidor. Isso inclui uma variedade de ferramentas, como Hudi e Iceberg, para trabalhar em grandes conjuntos de dados e usar Python e bibliotecas Python para enviar trabalhos do Spark.

Tópicos

- [Uso do Java 17 com o Amazon EMR Sem Servidor](#)
- [Uso do Apache Hudi com o EMR Sem Servidor](#)
- [Uso do Apache Iceberg com o EMR Sem Servidor](#)
- [Uso de bibliotecas Python com o EMR Sem Servidor](#)
- [Uso de diferentes versões do Python com o EMR Sem Servidor](#)
- [Usar o Delta Lake OSS com o EMR Sem Servidor](#)
- [Envio de trabalhos do EMR Sem Servidor provenientes do Airflow](#)
- [Uso de funções definidas pelo usuário no Hive com o EMR Sem Servidor](#)
- [Uso de imagens personalizadas com o EMR Sem Servidor](#)
- [Uso da integração do Amazon Redshift para Apache Spark no Amazon EMR Sem Servidor](#)
- [Como se conectar ao DynamoDB com o Amazon EMR Sem Servidor](#)

Uso do Java 17 com o Amazon EMR Sem Servidor

Com as versões 6.11.0 e posteriores do Amazon EMR, configure trabalhos do Spark no EMR Sem Servidor para usar o runtime do Java 17 na Java Virtual Machine (JVM). Use um dos métodos a seguir para configurar o Spark com Java 17.

JAVA_HOME

Para substituir a configuração da JVM para o EMR Sem Servidor 6.11.0 e superior, forneça a configuração `JAVA_HOME` para as classificações de ambiente `spark.emr-serverless.driverEnv` e `spark.executorEnv`.

x86_64

Defina as propriedades necessárias para especificar o Java 17 como a configuração `JAVA_HOME` para o driver e os executores do Spark:

```
--conf spark.emr-serverless.driverEnv.JAVA_HOME=/usr/lib/jvm/java-17-amazon-corretto.x86_64/
--conf spark.executorEnv.JAVA_HOME=/usr/lib/jvm/java-17-amazon-corretto.x86_64/
```

arm_64

Defina as propriedades necessárias para especificar o Java 17 como a configuração JAVA_HOME para o driver e os executores do Spark:

```
--conf spark.emr-serverless.driverEnv.JAVA_HOME=/usr/lib/jvm/java-17-amazon-corretto.aarch64/
--conf spark.executorEnv.JAVA_HOME=/usr/lib/jvm/java-17-amazon-corretto.aarch64/
```

spark-defaults

Como alternativa, você pode especificar o Java 17 na classificação spark-defaults para substituir a configuração da JVM no EMR Sem Servidor 6.11.0 e superior.

x86_64

Especifique Java 17 na classificação spark-defaults:

```
{
  "applicationConfiguration": [
    {
      "classification": "spark-defaults",
      "properties": {
        "spark.emr-serverless.driverEnv.JAVA_HOME" : "/usr/lib/jvm/java-17-amazon-corretto.x86_64/",
        "spark.executorEnv.JAVA_HOME": "/usr/lib/jvm/java-17-amazon-corretto.x86_64/"
      }
    }
  ]
}
```

arm_64

Especifique Java 17 na classificação spark-defaults:

```
{
```

```
"applicationConfiguration": [  
  {  
    "classification": "spark-defaults",  
    "properties": {  
      "spark.emr-serverless.driverEnv.JAVA_HOME" : "/usr/lib/jvm/java-17-  
amazon-corretto.aarch64/",  
      "spark.executorEnv.JAVA_HOME": "/usr/lib/jvm/java-17-amazon-  
corretto.aarch64/"  
    }  
  }  
]  
}
```

Uso do Apache Hudi com o EMR Sem Servidor

Esta seção descreve o uso do Apache Hudi com aplicações do EMR Sem Servidor. O Hudi é uma estrutura de gerenciamento de dados que simplifica o processamento de dados.

Para usar o Apache Hudi com aplicações do EMR Sem Servidor

1. Defina as propriedades necessárias do Spark na execução do trabalho correspondente do Spark.

```
spark.jars=/usr/lib/hudi/hudi-spark-bundle.jar,/usr/lib/hudi/hudi-utilities-  
bundle.jar,/usr/lib/hudi/hudi-aws-bundle.jar  
spark.serializer=org.apache.spark.serializer.KryoSerializer
```

2. Para sincronizar uma tabela Hudi com o catálogo configurado, designe o AWS Glue Data Catalog como sua metastore ou configure uma metastore externa. O EMR Sem Servidor é compatível com hms como modo de sincronização de tabelas do Hive para workloads do Hudi. O EMR Sem Servidor ativa essa propriedade como padrão. Para saber mais sobre como configurar a metastore, consulte [Configuração da metastore para EMR Sem Servidor](#).

Important

O EMR Sem Servidor não oferece suporte a HIVEQL ou JDBC como opções de modo de sincronização para tabelas do Hive lidarem com workloads do Hudi. Para saber mais, consulte [Modos de sincronização](#).

Ao usar o AWS Glue Data Catalog como sua metastore, especifique as seguintes propriedades de configuração para sua tarefa Hudi.

```
--conf spark.jars=/usr/lib/hudi/hudi-spark-bundle.jar,  
--conf spark.serializer=org.apache.spark.serializer.KryoSerializer,  
--conf  
spark.hadoop.hive.metastore.client.factory.class=com.amazonaws.glue.catalog.metastore.AWSG
```

Para saber mais sobre as versões de lançamento do Apache Hudi para o Amazon EMR, consulte [Histórico de lançamentos do Hudi](#).

Uso do Apache Iceberg com o EMR Sem Servidor

Esta seção descreve como usar o Apache Iceberg com aplicações do EMR Sem Servidor. O Apache Iceberg é um formato de tabela que ajuda a trabalhar com grandes conjuntos de dados em data lakes.

Para usar o Apache Iceberg com aplicações do EMR Sem Servidor

1. Defina as propriedades necessárias do Spark na execução do trabalho correspondente do Spark.

```
spark.jars=/usr/share/aws/iceberg/lib/iceberg-spark3-runtime.jar
```

2. Designe o AWS Glue Data Catalog como sua metastore ou configure uma metastore externa. Para saber mais sobre como configurar a metastore, consulte [Configuração da metastore para EMR Sem Servidor](#).

Configure as propriedades da metastore que você deseja usar no Iceberg. Por exemplo, se você quiser usar o AWS Glue Data Catalog, defina as seguintes propriedades na configuração do aplicativo.

```
spark.sql.catalog.dev.warehouse=s3://amzn-s3-demo-bucket/EXAMPLE-PREFIX/  
spark.sql.extensions=org.apache.iceberg.spark.extensions.IcebergSparkSessionExtensions  
spark.sql.catalog.dev=org.apache.iceberg.spark.SparkCatalog  
spark.sql.catalog.dev.catalog-impl=org.apache.iceberg.aws.glue.GlueCatalog  
spark.hadoop.hive.metastore.client.factory.class=com.amazonaws.glue.catalog.metastore.AWSG
```

Ao usar o AWS Glue Data Catalog como seu metastore, especifique as seguintes propriedades de configuração para seu trabalho no Iceberg.

```
--conf spark.jars=/usr/share/aws/iceberg/lib/iceberg-spark3-runtime.jar,  
--conf  
  spark.sql.extensions=org.apache.iceberg.spark.extensions.IcebergSparkSessionExtensions,  
--conf spark.sql.catalog.dev=org.apache.iceberg.spark.SparkCatalog,  
--conf spark.sql.catalog.dev.catalog-impl=org.apache.iceberg.aws.glue.GlueCatalog,  
--conf spark.sql.catalog.dev.warehouse=s3://amzn-s3-demo-bucket/EXAMPLE-PREFIX/  
--conf  
  spark.hadoop.hive.metastore.client.factory.class=com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalogHiveClient
```

Para saber mais sobre as versões do Apache Iceberg para o Amazon EMR, consulte [Histórico de lançamentos do Iceberg](#).

Uso de bibliotecas Python com o EMR Sem Servidor

Ao executar PySpark trabalhos em aplicativos sem servidor do Amazon EMR, empacote várias bibliotecas Python como dependências. Para fazer isso, use recursos nativos do Python, crie um ambiente virtual ou configure diretamente seus PySpark trabalhos para usar bibliotecas do Python. Esta página fala sobre cada abordagem.

Uso dos recursos nativos do Python

Ao definir a configuração a seguir, use PySpark para fazer upload de arquivos Python (.py), pacotes Python compactados () e arquivos Egg .zip () para os executores do Spark.egg.

```
--conf spark.submit.pyFiles=s3://amzn-s3-demo-bucket/EXAMPLE-PREFIX/<.py|.egg|.zip  
file>
```

Para obter mais detalhes sobre como usar ambientes virtuais Python para PySpark trabalhos, consulte [Usando recursos PySpark nativos](#).

Ao usar o Caderno do EMR, você pode disponibilizar a dependência do Python no seu Caderno executando o seguinte código:

```
%%configure -f
```

```
{
  "conf": {
    "spark.submit.pyFiles": "s3://amzn-s3-demo-bucket/EXAMPLE-PREFIX/<.py|.egg|.zip
file>
  }
}
```

Criação de um ambiente virtual Python

Para empacotar várias bibliotecas Python para um PySpark trabalho, crie ambientes virtuais Python isolados.

1. Para criar o ambiente virtual em Python, use os comandos a seguir. O exemplo mostrado instala os pacotes `scipy` e `matplotlib` em um pacote de ambiente virtual e copia o arquivo para um local do Amazon S3.

Important

Você deve executar os comandos a seguir em um ambiente do Amazon Linux 2 similar com a mesma versão do Python usada no EMR Sem Servidor, ou seja, Python 3.7.10 para o Amazon EMR versão 6.6.0. Você pode encontrar um exemplo de Dockerfile no repositório [EMR Serverless Samples](#). GitHub

```
# initialize a python virtual environment
python3 -m venv pyspark_venvsource
source pyspark_venvsource/bin/activate

# optionally, ensure pip is up-to-date
pip3 install --upgrade pip

# install the python packages
pip3 install scipy
pip3 install matplotlib

# package the virtual environment into an archive
pip3 install venv-pack
venv-pack -f -o pyspark_venv.tar.gz

# copy the archive to an S3 location
aws s3 cp pyspark_venv.tar.gz s3://amzn-s3-demo-bucket/EXAMPLE-PREFIX/
```

```
# optionally, remove the virtual environment directory
rm -fr pyspark_venvsource
```

2. Envie o trabalho do Spark com suas propriedades definidas para usar o ambiente virtual do Python.

```
--conf spark.archives=s3://amzn-s3-demo-bucket/EXAMPLE-PREFIX/
pyspark_venv.tar.gz#environment
--conf spark.emr-serverless.driverEnv.PYSPARK_DRIVER_PYTHON=./environment/bin/
python
--conf spark.emr-serverless.driverEnv.PYSPARK_PYTHON=./environment/bin/python
--conf spark.executorEnv.PYSPARK_PYTHON=./environment/bin/python
```

Observe que, se você não substituir o binário original do Python, a segunda configuração na sequência de configurações anterior será `--conf spark.executorEnv.PYSPARK_PYTHON=python`.

Para saber mais sobre como usar ambientes virtuais Python para PySpark trabalhos, consulte [Usando](#) o Virtualenv. Para obter mais exemplos de como enviar trabalhos do Spark, consulte [Uso das configurações do Spark ao executar trabalhos do EMR Sem Servidor](#).

Configurando PySpark trabalhos para usar bibliotecas Python

[Com as versões 6.12.0 e posteriores do Amazon EMR, você pode configurar diretamente trabalhos do EMR Serverless PySpark para usar bibliotecas Python populares de ciência de dados, como pandas, e sem nenhuma configuração adicional. NumPyPyArrow](#)

Os exemplos a seguir demonstram como empacotar cada biblioteca Python para um PySpark trabalho.

NumPy

NumPy é uma biblioteca Python para computação científica que oferece matrizes e operações multidimensionais para matemática, classificação, simulação aleatória e estatísticas básicas. Para usar NumPy, execute o seguinte comando:

```
import numpy
```

pandas

pandas é uma biblioteca Python baseada em NumPy. A biblioteca pandas fornece aos cientistas de dados estruturas de [DataFrames](#) e ferramentas de análise de dados. Para usar o pandas, execute o seguinte comando:

```
import pandas
```

PyArrow

PyArrow é uma biblioteca Python que gerencia dados colunares na memória para melhorar o desempenho no trabalho. PyArrow é baseado na especificação de desenvolvimento multilíngue Apache Arrow, que é uma forma padrão de representar e trocar dados em um formato colunar. Para usar PyArrow, execute o seguinte comando:

```
import pyarrow
```

Uso de diferentes versões do Python com o EMR Sem Servidor

Além do caso de uso em [Uso de bibliotecas Python com o EMR Sem Servidor](#), você também pode usar ambientes virtuais do Python para trabalhar com versões do Python diferentes da versão empacotada na versão do Amazon EMR para a aplicação do Amazon EMR. Para fazer isso, crie um ambiente virtual do Python com a versão do Python que deseja usar.

Para enviar um trabalho de um ambiente virtual do Python

1. Crie seu ambiente virtual com os comandos do exemplo a seguir. Este exemplo instala o Python 3.9.9 em um pacote de ambiente virtual e copia o arquivo para um local do Amazon S3.

Important

Se você usa o Amazon EMR nas versões 7.0.0 e superiores, execute seus comandos em um ambiente do Amazon Linux 2023 semelhante ao que você usa nas aplicações do EMR Sem Servidor.

Se você usa a versão 6.15.0 ou inferior, execute os comandos a seguir em um ambiente do Amazon Linux 2 similar.

```

# install Python 3.9.9 and activate the venv
yum install -y gcc openssl-devel bzip2-devel libffi-devel tar gzip wget make
wget https://www.python.org/ftp/python/3.9.9/Python-3.9.9.tgz && \
tar xzf Python-3.9.9.tgz && cd Python-3.9.9 && \
./configure --enable-optimizations && \
make altinstall

# create python venv with Python 3.9.9
python3.9 -m venv pyspark_venv_python_3.9.9 --copies
source pyspark_venv_python_3.9.9/bin/activate

# copy system python3 libraries to venv
cp -r /usr/local/lib/python3.9/* ./pyspark_venv_python_3.9.9/lib/python3.9/

# package venv to archive.
# **Note** that you have to supply --python-prefix option
# to make sure python starts with the path where your
# copied libraries are present.
# Copying the python binary to the "environment" directory.
pip3 install venv-pack
venv-pack -f -o pyspark_venv_python_3.9.9.tar.gz --python-prefix /home/hadoop/
environment

# stage the archive in S3
aws s3 cp pyspark_venv_python_3.9.9.tar.gz s3://<path>

# optionally, remove the virtual environment directory
rm -fr pyspark_venv_python_3.9.9

```

2. Defina suas propriedades para usar o ambiente virtual do Python e enviar o trabalho do Spark.

```

# note that the archive suffix "environment" is the same as the directory where you
# copied the Python binary.
--conf spark.archives=s3://amzn-s3-demo-bucket/EXAMPLE-PREFIX/
pyspark_venv_python_3.9.9.tar.gz#environment
--conf spark.emr-serverless.driverEnv.PYSPARK_DRIVER_PYTHON=./environment/bin/
python
--conf spark.emr-serverless.driverEnv.PYSPARK_PYTHON=./environment/bin/python
--conf spark.executorEnv.PYSPARK_PYTHON=./environment/bin/python

```

Para saber mais sobre como usar ambientes virtuais Python para PySpark trabalhos, consulte [Usando](#) o Virtualenv. Para obter mais exemplos de como enviar trabalhos do Spark, consulte [Uso das configurações do Spark ao executar trabalhos do EMR Sem Servidor](#).

Usar o Delta Lake OSS com o EMR Sem Servidor

Amazon EMR 6.9.0 e versões posteriores

Note

O Amazon EMR 7.0.0 e versões superiores usam o Delta Lake 3.0.0, que renomeia o arquivo `delta-core.jar` como `delta-spark.jar`. Se você usa o Amazon EMR 7.0.0 ou superior, certifique-se de especificar `delta-spark.jar` nas configurações.

O Amazon EMR 6.9.0 e versões posteriores incluem o Delta Lake, então você não precisa mais empacotar o Delta Lake por conta própria ou fornecer o sinalizador `--packages` com trabalhos do EMR Sem Servidor.

1. Ao enviar trabalhos do EMR Sem Servidor, certifique-se de ter as propriedades de configuração a seguir e incluir os parâmetros no campo `sparkSubmitParameters`.

```
--conf spark.jars=/usr/share/aws/delta/lib/delta-core.jar,/usr/share/aws/delta/lib/delta-storage.jar
--conf spark.sql.extensions=io.delta.sql.DeltaSparkSessionExtension
--conf
spark.sql.catalog.spark_catalog=org.apache.spark.sql.delta.catalog.DeltaCatalog
```

2. Crie um `delta_sample.py` local para testar a criação e a leitura de uma tabela do Delta.

```
# delta_sample.py
from pyspark.sql import SparkSession

import uuid

url = "s3://amzn-s3-demo-bucket/delta-lake/output/%s/" % str(uuid.uuid4())
spark = SparkSession.builder.appName("DeltaSample").getOrCreate()

## creates a Delta table and outputs to target S3 bucket
spark.range(5).write.format("delta").save(url)
```

```
## reads a Delta table and outputs to target S3 bucket
spark.read.format("delta").load(url).show
```

3. Usando o AWS CLI, faça o upload do `delta_sample.py` arquivo em seu bucket do Amazon S3. Em seguida, use o comando `start-job-run` para enviar um trabalho a uma aplicação existente do EMR Sem Servidor.

```
aws s3 cp delta_sample.py s3://amzn-s3-demo-bucket/code/

aws emr-serverless start-job-run \
  --application-id application-id \
  --execution-role-arn job-role-arn \
  --name emr-delta \
  --job-driver '{
    "sparkSubmit": {
      "entryPoint": "s3://amzn-s3-demo-bucket/code/delta_sample.py",
      "sparkSubmitParameters": "--conf spark.jars=/usr/share/
aws/delta/lib/delta-core.jar,/usr/share/aws/delta/lib/delta-storage.jar --
conf spark.sql.extensions=io.delta.sql.DeltaSparkSessionExtension --conf
spark.sql.catalog.spark_catalog=org.apache.spark.sql.delta.catalog.DeltaCatalog"
    }
  }'
```

Para usar bibliotecas do Python com o Delta Lake, adicione a biblioteca `delta-core` [empacotando-a como uma dependência](#) ou [usando-a como uma imagem personalizada](#).

Como alternativa, você pode usar `SparkContext.addPyFile` para adicionar as bibliotecas do Python do arquivo JAR `delta-core`:

```
import glob
from pyspark.sql import SparkSession

spark = SparkSession.builder.getOrCreate()
spark.sparkContext.addPyFile(glob.glob("/usr/share/aws/delta/lib/delta-core_*.jar")[0])
```

Amazon EMR 6.8.0 e versões anteriores

Se estiver usando o Amazon EMR 6.8.0 ou inferior, siga estas etapas para usar o OSS do Delta Lake com as aplicações do EMR Sem Servidor.

1. Para criar uma versão de código aberto do [Delta Lake](#) que seja compatível com a versão do Spark em seu aplicativo Amazon EMR Serverless, navegue até [o GitHub Delta e siga as instruções](#).
2. Faça o upload das bibliotecas do Delta Lake em um bucket do Amazon S3 em seu. Conta da AWS
3. Ao enviar trabalhos do EMR Sem Servidor na configuração da aplicação, inclua os arquivos JAR do Delta Lake que agora estão no bucket.

```
--conf spark.jars=s3://amzn-s3-demo-bucket/jars/delta-core_2.12-1.1.0.jar
```

4. Para garantir que você possa ler e gravar em uma tabela Delta, execute um PySpark teste de amostra.

```
from pyspark import SparkConf, SparkContext
    from pyspark.sql import HiveContext, SparkSession

import uuid

conf = SparkConf()
sc = SparkContext(conf=conf)
sqlContext = HiveContext(sc)

url = "s3://amzn-s3-demo-bucket/delta-lake/output/1.0.1/%s/" %
str(uuid.uuid4())

## creates a Delta table and outputs to target S3 bucket
session.range(5).write.format("delta").save(url)

## reads a Delta table and outputs to target S3 bucket
session.read.format("delta").load(url).show
```

Envio de trabalhos do EMR Sem Servidor provenientes do Airflow

O provedor da Amazon no Apache Airflow fornece operadores do EMR Sem Servidor. Para obter mais informações sobre operadores, consulte [Operadores do Amazon EMR Sem Servidor](#) na documentação do Apache Airflow.

Você pode usar `EmrServerlessCreateApplicationOperator` para criar uma aplicação do Spark ou do Hive. Você também pode usar `EmrServerlessStartJobOperator` para iniciar um ou mais trabalhos com a nova aplicação.

Para usar o operador com o Amazon Managed Workflows para Apache Airflow (MWAA) com Airflow 2.2.2, adicione a linha a seguir ao arquivo `requirements.txt` e atualize o ambiente do MWAA para usar o novo arquivo.

```
apache-airflow-providers-amazon==6.0.0
boto3>=1.23.9
```

Observe que o suporte ao EMR Sem Servidor foi adicionado à versão 5.0.0 do provedor da Amazon. A versão 6.0.0 é a última compatível com o Airflow 2.2.2. Você pode usar versões posteriores com o Airflow 2.4.3 no MWAA.

O exemplo abreviado a seguir mostra como criar uma aplicação, executar vários trabalhos do Spark e, em seguida, interromper a aplicação. Um exemplo completo está disponível no repositório [EMR Serverless Samples](#). GitHub Para obter detalhes adicionais sobre a configuração `sparkSubmit`, consulte [Uso das configurações do Spark ao executar trabalhos do EMR Sem Servidor](#).

```
from datetime import datetime

from airflow import DAG
from airflow.providers.amazon.aws.operators.emr import (
    EmrServerlessCreateApplicationOperator,
    EmrServerlessStartJobOperator,
    EmrServerlessDeleteApplicationOperator,
)

# Replace these with your correct values
JOB_ROLE_ARN = "arn:aws:iam::account-id:role/emr_serverless_default_role"
S3_LOGS_BUCKET = "amzn-s3-demo-bucket"

DEFAULT_MONITORING_CONFIG = {
    "monitoringConfiguration": {
        "s3MonitoringConfiguration": {"logUri": f"s3://amzn-s3-demo-bucket/logs/"}
    },
}

with DAG(
    dag_id="example_endtoend_emr_serverless_job",
    schedule_interval=None,
```

```

start_date=datetime(2021, 1, 1),
tags=["example"],
catchup=False,
) as dag:
    create_app = EmrServerlessCreateApplicationOperator(
        task_id="create_spark_app",
        job_type="SPARK",
        release_label="emr-6.7.0",
        config={"name": "airflow-test"},
    )

    application_id = create_app.output

    job1 = EmrServerlessStartJobOperator(
        task_id="start_job_1",
        application_id=application_id,
        execution_role_arn=JOB_ROLE_ARN,
        job_driver={
            "sparkSubmit": {
                "entryPoint": "local:///usr/lib/spark/examples/src/main/python/
pi_fail.py",
            }
        },
        configuration_overrides=DEFAULT_MONITORING_CONFIG,
    )

    job2 = EmrServerlessStartJobOperator(
        task_id="start_job_2",
        application_id=application_id,
        execution_role_arn=JOB_ROLE_ARN,
        job_driver={
            "sparkSubmit": {
                "entryPoint": "local:///usr/lib/spark/examples/src/main/python/pi.py",
                "entryPointArguments": ["1000"]
            }
        },
        configuration_overrides=DEFAULT_MONITORING_CONFIG,
    )

    delete_app = EmrServerlessDeleteApplicationOperator(
        task_id="delete_app",
        application_id=application_id,
        trigger_rule="all_done",
    )

```

```
(create_app >> [job1, job2] >> delete_app)
```

Uso de funções definidas pelo usuário no Hive com o EMR Sem Servidor

As funções definidas pelo usuário (UDFs) do Hive permitem criar funções personalizadas para processar registros ou grupos de registros. Neste tutorial, você usará um exemplo de UDF com uma aplicação do Amazon EMR Sem Servidor preexistente para executar um trabalho que gera um resultado de consulta. Para saber como configurar uma aplicação, consulte [Conceitos básicos do Amazon EMR Sem Servidor](#).

Para usar uma UDF com o EMR Sem Servidor

1. Navegue até o [GitHub](#) para obter uma amostra de UDF. Clone o repositório e alterne para a ramificação git que deseja usar. Atualize o maven-compiler-plugin no arquivo pom.xml do repositório para ter uma fonte. Atualize também a configuração da versão Java de destino para 1.8. Execute `mvn package -DskipTests` para criar o arquivo JAR que contém sua amostra UDFs.
2. Depois de criar o arquivo JAR, faça upload dele no bucket do S3 com o comando a seguir.

```
aws s3 cp brickhouse-0.8.2-JS.jar s3://amzn-s3-demo-bucket/jars/
```

3. Crie um exemplo de arquivo para usar um dos exemplos de funções de UDF. Salve a consulta como `udf_example.q` e faça upload dela no bucket do S3.

```
add jar s3://amzn-s3-demo-bucket/jars/brickhouse-0.8.2-JS.jar;
CREATE TEMPORARY FUNCTION from_json AS 'brickhouse.udf.json.FromJsonUDF';
select from_json('{"key1":[0,1,2], "key2":[3,4,5,6], "key3":[7,8,9]}', map("",
  array(cast(0 as int))));
select from_json('{"key1":[0,1,2], "key2":[3,4,5,6], "key3":[7,8,9]}', map("",
  array(cast(0 as int))))["key1"][2];
```

4. Envie o trabalho do Hive a seguir.

```
aws emr-serverless start-job-run \  
  --application-id application-id \  
  --execution-role-arn job-role-arn \  
  --job-driver '{
```

```

    "hive": {
      "query": "s3://amzn-s3-demo-bucket/queries/udf_example.q",
      "parameters": "--hiveconf hive.exec.scratchdir=s3://amzn-s3-demo-bucket/emr-serverless-hive/scratch --hiveconf hive.metastore.warehouse.dir=s3://'$BUCKET'/emr-serverless-hive/warehouse"
    }
  }' --configuration-overrides '{
    "applicationConfiguration": [{
      "classification": "hive-site",
      "properties": {
        "hive.driver.cores": "2",
        "hive.driver.memory": "6G"
      }
    }],
    "monitoringConfiguration": {
      "s3MonitoringConfiguration": {
        "logUri": "s3://amzn-s3-demo-bucket/logs/"
      }
    }
  }'

```

- Use o comando `get-job-run` para verificar o estado do trabalho. Espere até que o estado mude para SUCCESS.

```
aws emr-serverless get-job-run --application-id application-id --job-run-id job-id
```

- Faça download dos arquivos de saída com o comando a seguir.

```
aws s3 cp --recursive s3://amzn-s3-demo-bucket/logs/applications/application-id/jobs/job-id/HIVE_DRIVER/ .
```

O arquivo `stdout.gz` tem o aspecto a seguir.

```
{"key1": [0, 1, 2], "key2": [3, 4, 5, 6], "key3": [7, 8, 9]}
2
```

Uso de imagens personalizadas com o EMR Sem Servidor

Tópicos

- [Uso de uma versão personalizada do Python](#)

- [Uso de uma versão personalizada do Java](#)
- [Criação de uma imagem de ciência de dados](#)
- [Processamento de dados geoespaciais com o Apache Sedona](#)
- [Informações de licenciamento para o uso de imagens personalizadas](#)

Uso de uma versão personalizada do Python

Você pode criar uma imagem personalizada para usar uma versão diferente do Python. Para usar o Python versão 3.10 em trabalhos do Spark, por exemplo, execute o seguinte comando:

```
FROM public.ecr.aws/emr-serverless/spark/emr-6.9.0:latest

USER root

# install python 3
RUN yum install -y gcc openssl-devel bzip2-devel libffi-devel tar gzip wget make
RUN wget https://www.python.org/ftp/python/3.10.0/Python-3.10.0.tgz && \
tar xzf Python-3.10.0.tgz && cd Python-3.10.0 && \
./configure --enable-optimizations && \
make altinstall

# EMRS runs the image as hadoop
USER hadoop:hadoop
```

Antes de enviar o trabalho do Spark, defina as propriedades para usar o ambiente virtual do Python, da maneira mostrada a seguir.

```
--conf spark.emr-serverless.driverEnv.PYSPARK_DRIVER_PYTHON=/usr/local/bin/python3.10
--conf spark.emr-serverless.driverEnv.PYSPARK_PYTHON=/usr/local/bin/python3.10
--conf spark.executorEnv.PYSPARK_PYTHON=/usr/local/bin/python3.10
```

Uso de uma versão personalizada do Java

O exemplo a seguir demonstra como criar uma imagem personalizada para usar o Java 11 em trabalhos do Spark.

```
FROM public.ecr.aws/emr-serverless/spark/emr-6.9.0:latest
```

```
USER root

# install JDK 11
RUN amazon-linux-extras install java-openjdk11

# EMRS runs the image as hadoop
USER hadoop:hadoop
```

Antes de enviar o trabalho do Spark, defina as propriedades do Spark para usar o Java 11, da maneira mostrada a seguir.

```
--conf spark.executorEnv.JAVA_HOME=/usr/lib/jvm/java-11-
openjdk-11.0.16.0.8-1.amzn2.0.1.x86_64
--conf spark.emr-serverless.driverEnv.JAVA_HOME=/usr/lib/jvm/java-11-
openjdk-11.0.16.0.8-
```

Criação de uma imagem de ciência de dados

O exemplo a seguir mostra como incluir pacotes Python comuns de ciência de dados, como Pandas e NumPy

```
FROM public.ecr.aws/emr-serverless/spark/emr-6.9.0:latest

USER root

# python packages
RUN pip3 install boto3 pandas numpy
RUN pip3 install -U scikit-learn==0.23.2 scipy
RUN pip3 install sk-dist
RUN pip3 install xgboost

# EMR Serverless runs the image as hadoop
USER hadoop:hadoop
```

Processamento de dados geoespaciais com o Apache Sedona

O exemplo a seguir mostra como criar uma imagem para incluir o Apache Sedona para processamento geoespacial.

```
FROM public.ecr.aws/emr-serverless/spark/emr-6.9.0:latest
```

```
USER root

RUN yum install -y wget
RUN wget https://repo1.maven.org/maven2/org/apache/sedona/sedona-core-3.0_2.12/1.3.0-incubating/sedona-core-3.0_2.12-1.3.0-incubating.jar -P /usr/lib/spark/jars/
RUN pip3 install apache-sedona

# EMRS runs the image as hadoop
USER hadoop:hadoop
```

Informações de licenciamento para o uso de imagens personalizadas

Você pode criar imagens personalizadas com o EMR Sem Servidor para realizar tarefas específicas ou usar versões específicas de um pacote de software. A modificação e a distribuição de imagens personalizadas podem estar sujeitas a regras e termos de licenciamento. O texto de licenciamento aparece na subseção abaixo.

Licenciamento aplicável a imagens personalizadas

Copyright Amazon.com e suas afiliadas; todos os direitos reservados. Este software é um AWS Conteúdo de acordo com o [Contrato AWS do Cliente](#) e não pode ser distribuído sem permissão. Além das permissões na [Licença de AWS Propriedade Intelectual](#), o AWS Licenciador concede a você estas permissões adicionais:

É permitido criar, copiar e usar derivados do AWS conteúdo, desde que as seguintes condições sejam atendidas:

- Você não modifica o AWS conteúdo em si, e quaisquer derivados são estritamente o resultado da adição de novos conteúdos por você.
- Reproduções internas devem manter o aviso de direitos autorais acima.
- A distribuição externa, em formato fonte ou binário, com ou sem modificações, não é permitida nos termos desta licença.

Para obter mais informações sobre o uso de imagens personalizadas, consulte [Usar imagens personalizadas com o EMR Sem Servidor](#).

Uso da integração do Amazon Redshift para Apache Spark no Amazon EMR Sem Servidor

Com as versões 6.9.0 e posteriores do Amazon EMR, cada imagem de versão inclui um conector entre o [Apache Spark](#) e o Amazon Redshift. Com esse conector, use o Spark no Amazon EMR Sem Servidor para processar dados armazenados no Amazon Redshift. A integração é baseada no [conector de código aberto spark-redshift](#). No Amazon EMR Sem Servidor, a [integração do Amazon Redshift para Apache Spark](#) está incluída como uma integração nativa.

Tópicos

- [Inicialização de uma aplicação do Spark com a integração do Amazon Redshift para Apache Spark](#)
- [Autenticação com a integração do Amazon Redshift para Apache Spark](#)
- [Leitura e gravação de e para o Amazon Redshift](#)
- [Considerações e limitações ao usar o conector do Spark](#)

Inicialização de uma aplicação do Spark com a integração do Amazon Redshift para Apache Spark

Para usar a integração com o EMR Sem Servidor 6.9.0, passe as dependências necessárias do Redshift para Spark com o trabalho do Spark. Use `--jars` para incluir as bibliotecas relacionadas ao conector do Redshift. Para acessar outros locais de arquivo com suporte pela opção `--jars`, consulte a seção [Gerenciamento avançado de dependências](#) na documentação do Apache Spark.

- `spark-redshift.jar`
- `spark-avro.jar`
- `RedshiftJDBC.jar`
- `minimal-json.jar`

As versões 6.10.0 e superiores do Amazon EMR não exigem a dependência `minimal-json.jar` e, por padrão, instalam automaticamente as outras dependências em cada cluster. Os exemplos a seguir demonstram como iniciar uma aplicação do Spark com a integração do Amazon Redshift para Apache Spark.

Amazon EMR 6.10.0 +

Inicie um trabalho do Spark no Amazon EMR Sem Servidor com a integração do Amazon Redshift para Apache Spark na versão 6.10.0 e posteriores do Amazon EMR Sem Servidor.

```
spark-submit my_script.py
```

Amazon EMR 6.9.0

Para executar um trabalho do Spark no Amazon EMR Sem Servidor com a integração do Amazon Redshift para Apache Spark no EMR Sem Servidor versão 6.9.0, use a opção `--jars` como mostrado no exemplo a seguir. Observe que os caminhos listados com a opção `--jars` são os caminhos padrão para os arquivos JAR.

```
--jars
  /usr/share/aws/redshift/jdbc/RedshiftJDBC.jar,
  /usr/share/aws/redshift/spark-redshift/lib/spark-redshift.jar,
  /usr/share/aws/redshift/spark-redshift/lib/spark-avro.jar,
  /usr/share/aws/redshift/spark-redshift/lib/minimal-json.jar
```

```
spark-submit \
  --jars /usr/share/aws/redshift/jdbc/RedshiftJDBC.jar,/usr/share/aws/redshift/
spark-redshift/lib/spark-redshift.jar,/usr/share/aws/redshift/spark-redshift/lib/
spark-avro.jar,/usr/share/aws/redshift/spark-redshift/lib/minimal-json.jar \
  my_script.py
```

Autenticação com a integração do Amazon Redshift para Apache Spark

Use AWS Secrets Manager para recuperar credenciais e conectar-se ao Amazon Redshift

Você pode se autenticar com segurança no Amazon Redshift armazenando as credenciais no Secrets Manager e fazer com que o trabalho do Spark chame a API `GetSecretValue` para buscá-la:

```
from pyspark.sql import SQLContext
import boto3

sc = # existing SparkContext
sql_context = SQLContext(sc)
```

```
secretsmanager_client = boto3.client('secretsmanager',
    region_name=os.getenv('AWS_REGION'))
secret_manager_response = secretsmanager_client.get_secret_value(
    SecretId='string',
    VersionId='string',
    VersionStage='string'
)
username = # get username from secret_manager_response
password = # get password from secret_manager_response
url = "jdbc:redshift://redshifthost:5439/database?user=" + username + "&password="
    + password

# Access to Redshift cluster using Spark
```

Autenticação no Amazon Redshift com um driver JDBC

Definição de um nome de usuário e de uma senha no URL do JDBC

Você pode autenticar um trabalho do Spark em um cluster do Amazon Redshift especificando o nome e a senha do banco de dados do Amazon Redshift no URL do JDBC.

Note

Se você transferir as credenciais do banco de dados no URL, qualquer pessoa que tenha acesso ao URL também poderá acessar as credenciais. Este método geralmente não é recomendado porque não é uma opção segura.

Se a segurança não for uma preocupação para sua aplicação, use o seguinte formato para definir o nome de usuário e a senha no URL do JDBC:

```
jdbc:redshift://redshifthost:5439/database?user=username&password=password
```

Uso da autenticação baseada no IAM com o perfil de execução de trabalho do Amazon EMR Sem Servidor

A partir da versão 6.9.0 do Amazon EMR Sem Servidor, a versão 2.1 ou superior do driver JDBC do Amazon Redshift é empacotada no ambiente. Com a versão 2.1 e versões superiores do driver JDBC, é possível especificar o URL do JDBC e não incluir o nome de usuário e a senha brutos.

Em vez disso, especifique o esquema `jdbc:redshift:iam://`. Isso comanda o driver JDBC para usar o perfil de execução de trabalho do EMR Sem Servidor para buscar as credenciais automaticamente. Consulte [Configurar uma conexão JDBC ou ODBC para usar credenciais do IAM](#) no Guia de gerenciamento do Amazon Redshift para obter mais informações. Um exemplo desse URL é:

```
jdbc:redshift:iam://examplecluster.abc123xyz789.us-west-2.redshift.amazonaws.com:5439/dev
```

As seguintes permissões são obrigatórias para o perfil de execução do trabalho quando ele atende às condições fornecidas:

Permissão	Condições para se tornar obrigatória para o perfil de execução de trabalho
<code>redshift:GetClusterCredentials</code>	Obrigatória para que o driver JDBC busque as credenciais do Amazon Redshift.
<code>redshift:DescribeCluster</code>	Obrigatória se você especificar o cluster do Amazon Redshift e a Região da AWS no URL do JDBC em vez do endpoint.
<code>redshift-serverless:GetCredentials</code>	Obrigatória para que o driver JDBC busque as credenciais do Amazon Redshift sem servidor.
<code>redshift-serverless:GetWorkgroup</code>	Obrigatória se você estiver usando o Amazon Redshift sem servidor e especificando o URL em termos de nome e de região do grupo de trabalho

Como se conectar ao Amazon Redshift em uma VPC diferente

Ao definir um cluster provisionado do Amazon Redshift ou um grupo de trabalho do Amazon Redshift Sem Servidor em uma VPC, configure a conectividade da VPC para que a aplicação do Amazon EMR Sem Servidor acesse os recursos. Para obter mais informações sobre como configurar a conectividade da VPC em uma aplicação do EMR Sem Servidor, consulte [Configuração do acesso à VPC para que aplicações do EMR Sem Servidor se conectem aos dados](#).

- Se o cluster provisionado do Amazon Redshift ou o grupo de trabalho do Amazon Redshift Sem Servidor estiver acessível publicamente, você poderá especificar uma ou mais sub-redes privadas que tenham um gateway NAT anexado ao criar aplicações do EMR Sem Servidor.
- Se o cluster provisionado do Amazon Redshift ou o grupo de trabalho do Amazon Redshift Sem Servidor não estiver acessível publicamente, você deverá criar um endpoint da VPC gerenciado pelo Amazon Redshift para o cluster do Amazon Redshift, conforme descrito em [Configuração do acesso à VPC para que aplicações do EMR Sem Servidor se conectem aos dados](#). Como alternativa, você pode criar o grupo de trabalho do Amazon Redshift Sem Servidor conforme descrito em [Conectar-se ao Amazon Redshift Serverless](#) no Guia de gerenciamento do Amazon Redshift. É necessário associar o cluster ou subgrupo às sub-redes privadas que você especifica ao criar a aplicação do EMR Sem Servidor.

Note

Se você usa autenticação baseada no IAM e suas sub-redes privadas na aplicação do EMR Sem Servidor não têm um gateway NAT anexado, também é necessário criar um endpoint da VPC nessas sub-redes para o Amazon Redshift ou o Amazon Redshift Sem Servidor. Dessa forma, o driver JDBC pode buscar as credenciais.

Leitura e gravação de e para o Amazon Redshift

Os exemplos de código a seguir são usados PySpark para ler e gravar dados de amostra de e para um banco de dados do Amazon Redshift com uma API de fonte de dados e com o SparkSQL.

Data source API

Use PySpark para ler e gravar dados de amostra de e para um banco de dados do Amazon Redshift com a API de fonte de dados.

```
import boto3
from pyspark.sql import SQLContext

sc = # existing SparkContext
sql_context = SQLContext(sc)

url = "jdbc:redshift:iam://redshifthost:5439/database"
aws_iam_role_arn = "arn:aws:iam::account-id:role/role-name"
```

```
df = sql_context.read \
    .format("io.github.spark_redshift_community.spark.redshift") \
    .option("url", url) \
    .option("dbtable", "table-name") \
    .option("tempdir", "s3://path/for/temp/data") \
    .option("aws_iam_role", "aws-iam-role-arn") \
    .load()

df.write \
    .format("io.github.spark_redshift_community.spark.redshift") \
    .option("url", url) \
    .option("dbtable", "table-name-copy") \
    .option("tempdir", "s3://path/for/temp/data") \
    .option("aws_iam_role", "aws-iam-role-arn") \
    .mode("error") \
    .save()
```

SparkSQL

Use PySpark para ler e gravar dados de amostra de e para um banco de dados do Amazon Redshift com o SparkSQL.

```
import boto3
import json
import sys
import os
from pyspark.sql import SparkSession

spark = SparkSession \
    .builder \
    .enableHiveSupport() \
    .getOrCreate()

url = "jdbc:redshift:iam://redshifthost:5439/database"
aws_iam_role_arn = "arn:aws:iam::account-id:role/role-name"

bucket = "s3://path/for/temp/data"
tableName = "table-name" # Redshift table name

s = f"""CREATE TABLE IF NOT EXISTS {table-name} (country string, data string)
USING io.github.spark_redshift_community.spark.redshift
OPTIONS (dbtable '{table-name}', tempdir '{bucket}', url '{url}', aws_iam_role
'{aws-iam-role-arn}'); """
```

```
spark.sql(s)

columns = ["country" ,"data"]
data = [("test-country", "test-data")]
df = spark.sparkContext.parallelize(data).toDF(columns)

# Insert data into table
df.write.insertInto(table-name, overwrite=False)
df = spark.sql(f"SELECT * FROM {table-name}")
df.show()
```

Considerações e limitações ao usar o conector do Spark

- Sugerimos que você ative o SSL para a conexão JDBC do Spark no Amazon EMR ao Amazon Redshift.
- Sugerimos que você gerencie as credenciais do cluster do Amazon Redshift no AWS Secrets Manager como uma prática recomendada. Consulte [Usando AWS Secrets Manager para recuperar credenciais para se conectar ao Amazon Redshift para ver](#) um exemplo.
- Sugerimos que você transmita um perfil do IAM com o parâmetro `aws_iam_role` para o parâmetro de autenticação do Amazon Redshift.
- No momento, o parâmetro `tempformat` não é compatível com o formato Parquet.
- O URI `tempdir` aponta para um local do Amazon S3. Esse diretório temporário não é limpo automaticamente e, portanto, pode incorrer em custos adicionais.
- Considere as seguintes recomendações para o Amazon Redshift:
 - Sugerimos bloquear o acesso público ao cluster do Amazon Redshift.
 - Sugerimos ativar o [registro em log de auditoria do Amazon Redshift](#).
 - Sugerimos que você ative a [Criptografia em repouso do Amazon Redshift](#).
- Considere as seguintes recomendações para o Amazon S3:
 - Sugerimos que você [bloqueie o acesso público aos buckets do Amazon S3](#).
 - Sugerimos que você use [criptografia no lado do servidor do Amazon S3](#) para criptografar os buckets do Amazon S3 usados.
 - Sugerimos que você use as [políticas de ciclo de vida do Amazon S3](#) para definir as regras de retenção para o bucket do Amazon S3.

- O Amazon EMR sempre verifica o código importado do código aberto para a imagem. Por motivos de segurança, não oferecemos suporte aos seguintes métodos de autenticação do Spark para o Amazon S3:
 - Definindo chaves de AWS acesso na classificação `hadoop-env` de configuração
 - Codificação de chaves de AWS acesso no URI `tempdir`

Para obter mais informações sobre como usar o conector e os parâmetros compatíveis, consulte os seguintes recursos:

- [Integração do Amazon Redshift para Apache Spark](#) no Guia de gerenciamento do Amazon Redshift.
- O [repositório da comunidade spark-redshift](#) no GitHub.

Como se conectar ao DynamoDB com o Amazon EMR Sem Servidor

Neste tutorial, você faz upload de um subconjunto de dados do [United States Board on Geographic Names](#) para um bucket do Amazon S3 e, em seguida, usa o Hive ou o Spark no Amazon EMR Sem Servidor para copiar os dados em uma tabela do Amazon DynamoDB para consulta.

Etapa 1: upload dos dados em um bucket do Amazon S3

Para criar um bucket do Amazon S3, siga as instruções em [Criação de um bucket](#) no Guia do usuário do console do Amazon Simple Storage Service. Substitua as referências a `amzn-s3-demo-bucket` pelo nome do bucket recém-criado. Agora, a aplicação do EMR Sem Servidor está pronta para executar trabalhos.

1. Faça download do arquivo de exemplo de dados `features.zip` com o comando a seguir.

```
wget https://docs.aws.amazon.com/amazondynamodb/latest/developerguide/samples/features.zip
```

2. Extraia o arquivo `features.txt` do arquivamento e acessar as primeiras linhas do arquivo:

```
unzip features.zip  
head features.txt
```

O resultado deve ser semelhante ao mostrado a seguir.

```
1535908|Big Run|Stream|WV|38.6370428|-80.8595469|794
875609|Constable Hook|Cape|NJ|40.657881|-74.0990309|7
1217998|Gooseberry Island|Island|RI|41.4534361|-71.3253284|10
26603|Boone Moore Spring|Spring|AZ|34.0895692|-111.410065|3681
1506738|Missouri Flat|Flat|WA|46.7634987|-117.0346113|2605
1181348|Minnow Run|Stream|PA|40.0820178|-79.3800349|1558
1288759|Hunting Creek|Stream|TN|36.343969|-83.8029682|1024
533060|Big Charles Bayou|Bay|LA|29.6046517|-91.9828654|0
829689|Greenwood Creek|Stream|NE|41.596086|-103.0499296|3671
541692|Button Willow Island|Island|LA|31.9579389|-93.0648847|98
```

Os campos em cada linha aqui indicam um identificador exclusivo, nome, tipo de característica natural, estado, latitude em graus, longitude em graus e altura em pés.

3. Upload de dados no Amazon S3

```
aws s3 cp features.txt s3://amzn-s3-demo-bucket/features/
```

Etapa 2: criar uma tabela do Hive

Use o Apache Spark ou o Hive para criar uma tabela do Hive que contenha os dados carregados no Amazon S3.

Spark

Para criar uma tabela do Hive com o Spark, execute o comando a seguir.

```
import org.apache.spark.sql.SparkSession

val sparkSession = SparkSession.builder().enableHiveSupport().getOrCreate()

sparkSession.sql("CREATE TABLE hive_features \
  (feature_id BIGINT, \
  feature_name STRING, \
  feature_class STRING, \
  state_alpha STRING, \
  prim_lat_dec DOUBLE, \
  prim_long_dec DOUBLE, \
```

```
elev_in_ft BIGINT) \  
ROW FORMAT DELIMITED \  
FIELDS TERMINATED BY '|' \  
LINES TERMINATED BY '\n' \  
LOCATION 's3://amzn-s3-demo-bucket/features';")
```

Agora você tem uma tabela do Hive preenchida com os dados do arquivo `features.txt`. Para verificar se seus dados estão na tabela, execute uma consulta do Spark SQL conforme mostrado no exemplo a seguir.

```
sparkSession.sql(  
    "SELECT state_alpha, COUNT(*) FROM hive_features GROUP BY state_alpha;")
```

Hive

Para criar uma tabela do Hive com o Hive, execute o comando a seguir.

```
CREATE TABLE hive_features  
    (feature_id          BIGINT,  
     feature_name        STRING ,  
     feature_class       STRING ,  
     state_alpha         STRING,  
     prim_lat_dec        DOUBLE ,  
     prim_long_dec       DOUBLE ,  
     elev_in_ft          BIGINT)  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY '|'\  
LINES TERMINATED BY '\n'  
LOCATION 's3://amzn-s3-demo-bucket/features';
```

Agora você tem uma tabela do Hive que contém dados do arquivo `features.txt`. Para verificar se seus dados estão na tabela, execute uma consulta do HiveQL, conforme mostrado no exemplo a seguir.

```
SELECT state_alpha, COUNT(*) FROM hive_features GROUP BY state_alpha;
```

Etapa 3: copiar dados para o DynamoDB

Use o Spark ou o Hive para copiar dados para uma nova tabela do DynamoDB.

Spark

Para copiar dados da tabela do Hive criada na etapa anterior para o DynamoDB, siga as etapas de 1 a 3 em [Copiar dados para o DynamoDB](#). Isso cria uma tabela do DynamoDB chamada Features. Em seguida, você pode ler os dados diretamente do arquivo de texto e copiá-los para a tabela do DynamoDB, conforme mostra o exemplo a seguir.

```
import com.amazonaws.services.dynamodbv2.model.AttributeValue
import org.apache.hadoop.dynamodb.DynamoDBItemWritable
import org.apache.hadoop.dynamodb.read.DynamoDBInputFormat
import org.apache.hadoop.io.Text
import org.apache.hadoop.mapred.JobConf
import org.apache.spark.SparkContext

import scala.collection.JavaConverters._

object EmrServerlessDynamoDbTest {

  def main(args: Array[String]): Unit = {

    jobConf.set("dynamodb.input.tableName", "Features")
    jobConf.set("dynamodb.output.tableName", "Features")
    jobConf.set("dynamodb.region", "region")

    jobConf.set("mapred.output.format.class",
"org.apache.hadoop.dynamodb.write.DynamoDBOutputFormat")
    jobConf.set("mapred.input.format.class",
"org.apache.hadoop.dynamodb.read.DynamoDBInputFormat")

    val rdd = sc.textFile("s3://amzn-s3-demo-bucket/ddb-connector/")
      .map(row => {
        val line = row.split("\\|")
        val item = new DynamoDBItemWritable()

        val elevInFt = if (line.length > 6) {
          new AttributeValue().withN(line(6))
        } else {
          new AttributeValue().withNULL(true)
        }

        item.setItem(Map(
          "feature_id" -> new AttributeValue().withN(line(0)),
          "feature_name" -> new AttributeValue(line(1)),
```

```

        "feature_class" -> new AttributeValue(line(2)),
        "state_alpha" -> new AttributeValue(line(3)),
        "prim_lat_dec" -> new AttributeValue().withN(line(4)),
        "prim_long_dec" -> new AttributeValue().withN(line(5)),
        "elev_in_ft" -> elevInFt)
        .asJava)
    (new Text(""), item)
    })
    rdd.saveAsHadoopDataset(jobConf)
}
}

```

Hive

Para copiar dados da tabela do Hive criados na etapa anterior para o DynamoDB, siga as instruções em [Copiar dados para o DynamoDB](#).

Etapa 4: consultar dados do DynamoDB

Use o Spark ou o Hive para consultar sua tabela do DynamoDB.

Spark

Para consultar dados da tabela do DynamoDB que você criou na etapa anterior, use o Spark SQL ou a API Spark. MapReduce

Example— Consulte sua tabela do DynamoDB com o Spark SQL

A consulta do Spark SQL a seguir retorna uma lista de todos os tipos de recursos em ordem alfabética.

```

val dataframe = sparkSession.sql("SELECT DISTINCT feature_class \
FROM ddb_features \
ORDER BY feature_class;")

```

A consulta do Spark SQL a seguir retorna uma lista de todos os lakes que começam com a letra M.

```

val dataframe = sparkSession.sql("SELECT feature_name, state_alpha \
FROM ddb_features \
WHERE feature_class = 'Lake' \
AND feature_name LIKE 'M%' \

```

```
ORDER BY feature_name;")
```

A consulta do Spark SQL a seguir retorna uma lista de todos os estados com pelo menos três recursos que ultrapassam uma milha.

```
val dataframe = sparkSession.dql("SELECT state_alpha, feature_class, COUNT(*) \
FROM ddb_features \
WHERE elev_in_ft > 5280 \
GROUP BY state_alpha, feature_class \
HAVING COUNT(*) >= 3 \
ORDER BY state_alpha, feature_class;")
```

Example— Consulte sua tabela do DynamoDB com a API Spark MapReduce

A MapReduce consulta a seguir retorna uma lista de todos os tipos de recursos em ordem alfabética.

```
val df = sc.hadoopRDD(jobConf, classOf[DynamoDBInputFormat], classOf[Text],
  classOf[DynamoDBItemWritable])
  .map(pair => (pair._1, pair._2.getItem))
  .map(pair => pair._2.get("feature_class").getS)
  .distinct()
  .sortBy(value => value)
  .toDF("feature_class")
```

A MapReduce consulta a seguir retorna uma lista de todos os lagos que começam com a letra M.

```
val df = sc.hadoopRDD(jobConf, classOf[DynamoDBInputFormat], classOf[Text],
  classOf[DynamoDBItemWritable])
  .map(pair => (pair._1, pair._2.getItem))
  .filter(pair => "Lake".equals(pair._2.get("feature_class").getS))
  .filter(pair => pair._2.get("feature_name").getS.startsWith("M"))
  .map(pair => (pair._2.get("feature_name").getS,
  pair._2.get("state_alpha").getS))
  .sortBy(_._1)
  .toDF("feature_name", "state_alpha")
```

A MapReduce consulta a seguir retorna uma lista de todos os estados com pelo menos três características com mais de uma milha.

```
val df = sc.hadoopRDD(jobConf, classOf[DynamoDBInputFormat], classOf[Text],
  classOf[DynamoDBItemWritable])
```

```
.map(pair => pair._2.getItem)
.filter(pair => pair.get("elev_in_ft").getN != null)
.filter(pair => Integer.parseInt(pair.get("elev_in_ft").getN) > 5280)
.groupBy(pair => (pair.get("state_alpha").getS, pair.get("feature_class").getS))
.filter(pair => pair._2.size >= 3)
.map(pair => (pair._1._1, pair._1._2, pair._2.size))
.sortBy(pair => (pair._1, pair._2))
.toDF("state_alpha", "feature_class", "count")
```

Hive

Para consultar dados da tabela do DynamoDB criada na etapa anterior, siga as instruções em [Query the data in the DynamoDB table](#).

Configurar o acesso entre contas

Para configurar o acesso entre contas do EMR Sem Servidor, conclua as etapas a seguir. No exemplo, AccountA é a conta na qual você criou a aplicação do Amazon EMR Sem Servidor e AccountB é a conta em que o Amazon DynamoDB está localizado.

1. Crie uma tabela do DynamoDB em AccountB. Para obter mais informações, consulte [Etapa 1: criar uma tabela](#).
2. Crie um perfil do IAM Cross-Account-Role-B na AccountB que possa acessar a tabela do DynamoDB.
 - a. Faça login no Console de gerenciamento da AWS e abra o console do IAM em <https://console.aws.amazon.com/iam/>.
 - b. Escolha Perfis e crie um perfil chamado Cross-Account-Role-B. Para obter mais informações sobre como criar perfis do IAM, consulte [Criar perfis do IAM](#) no Guia do usuário do IAM.
 - c. Crie uma política do IAM que conceda permissões de acesso à tabela do DynamoDB entre contas. Em seguida, anexe a política do IAM ao Cross-Account-Role-B.

A seguir está uma política que concede acesso a uma tabela do DynamoDB CrossAccountTable.

- d. Edite a relação de confiança para o perfil Cross-Account-Role-B.

Para configurar a relação de confiança para a função, escolha a guia Relações de confiança no console do IAM para a função que você criou na Etapa 2: Cross-Account-Role-B.

Selecione Editar relacionamento de confiança e adicione o documento de política a seguir. Esse documento permite que Job-Execution-Role-A em AccountA assumo o perfil Cross-Account-Role-B.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "sts:AssumeRole"
      ],
      "Resource": "arn:aws:iam::123456789012:role/Job-Execution-Role-A",
      "Sid": "AllowSTSAssumerole"
    }
  ]
}
```

- e. Conceda a Job-Execution-Role-A em AccountA com - STS Assume role permissões para assumir Cross-Account-Role-B.

No console do IAM para Conta da AWS AccountA, selecione Job-Execution-Role-A. Adicione a instrução de política a seguir ao Job-Execution-Role-A para permitir a ação AssumeRole no perfil Cross-Account-Role-B.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "sts:AssumeRole"
      ],
      "Resource": [
        "arn:aws:iam::123456789012:role/Cross-Account-Role-B"
      ]
    }
  ]
}
```

```
        "Sid": "AllowSTSAssumerole"  
    }  
  ]  
}
```

- f. Defina a propriedade `dynamodb.customAWSCredentialsProvider` com valor como `com.amazonaws.emr.AssumeRoleAWSCredentialsProvider` na classificação do site principal. Defina a variável de ambiente `ASSUME_ROLE_CREDENTIALS_ROLE_ARN` com o valor do ARN de `Cross-Account-Role-B`.
3. Execute o trabalho do Spark ou do Hive usando `Job-Execution-Role-A`.

Considerações

Observe esses comportamentos e limitações ao usar o conector do DynamoDB com o Apache Spark ou o Apache Hive.

Considerações ao usar o conector do DynamoDB com o Apache Spark

- O Spark SQL não é compatível com a criação de uma tabela do Hive com a opção de manipulador de armazenamento. Para obter mais informações, consulte [Especificando o formato de armazenagem para tabelas Hive](#) na documentação do Apache Spark.
- O Spark SQL não é compatível com a operação `STORED BY` com o manipulador de armazenamento. Se você quiser interagir com uma tabela do DynamoDB por meio de uma tabela externa do Hive, use o Hive para criar a tabela primeiro.
- Para traduzir uma consulta em uma consulta do DynamoDB, o conector do DynamoDB usa o pushdown de predicado. O pushdown de predicado filtra os dados por uma coluna que é mapeada para a chave de partição de uma tabela do DynamoDB. O predicate pushdown só funciona quando você usa o conector com o Spark SQL, e não com a API. MapReduce

Considerações ao usar o conector do DynamoDB com o Apache Hive

Ajuste do número máximo de mapeadores

- Se você usar a consulta `SELECT` para ler dados de uma tabela externa do Hive mapeada no DynamoDB, o número de tarefas de mapeamento no EMR Sem Servidor será calculado como o throughput total de leitura configurado na tabela do DynamoDB, dividido pelo throughput por tarefa de mapeamento. O throughput padrão por tarefa de mapeamento é 100.

- O trabalho do Hive pode usar o número de tarefas de mapeamento além do número máximo de contêineres configurados por aplicação do EMR Sem Servidor, dependendo do throughput de leitura configurado para o DynamoDB. Além disso, uma consulta do Hive de longa duração pode consumir toda a capacidade de leitura provisionada da tabela do DynamoDB. Isso afeta negativamente outros usuários.
- Você pode usar a propriedade `dynamodb.max.map.tasks` para definir um limite superior para mapear tarefas. Você também pode usar essa propriedade para ajustar a quantidade de dados lidos por cada tarefa de mapa com base no tamanho do contêiner da tarefa.
- Você pode definir a propriedade `dynamodb.max.map.tasks` no nível de consulta do Hive ou na classificação `hive-site` do comando `start-job-run`. Esse valor deve ser igual ou maior que 1. Quando o Hive processa a consulta, o trabalho do Hive resultante não usa mais do que os valores de `dynamodb.max.map.tasks` ao ler a tabela do DynamoDB.

Ajuste do throughput de gravação por tarefa

- O throughput de gravação por tarefa no EMR Sem Servidor é calculado como o throughput total de gravação configurado para uma tabela do DynamoDB, dividido pelo valor da propriedade `mapreduce.job.maps`. No Hive, o valor padrão dessa propriedade é 2. Assim, as duas primeiras tarefas no estágio final do trabalho do Hive podem consumir todo o throughput de gravação. Isso leva ao controle de utilização das gravações de outras tarefas no mesmo trabalho ou em outros.
- Para evitar o controle de utilização de gravação, defina o valor da propriedade `mapreduce.job.maps` com base no número de tarefas no estágio final ou no throughput de gravação que você deseja alocar por tarefa. Defina essa propriedade na classificação `mapred-site` do comando `start-job-run` no EMR Sem Servidor.

Segurança

A segurança na nuvem AWS é a maior prioridade. Como AWS cliente, você se beneficia de uma arquitetura de data center e rede criada para atender aos requisitos das organizações mais sensíveis à segurança.

A segurança é uma responsabilidade compartilhada entre você AWS e você. O [modelo de responsabilidade compartilhada](#) descreve isto como segurança da nuvem e segurança na nuvem.

- Segurança da nuvem — AWS é responsável por proteger a infraestrutura que executa AWS os serviços na AWS nuvem. AWS também fornece serviços que você usa com segurança. Auditores de terceiros testam e verificam regularmente a eficácia da nossa segurança como parte dos [programas de conformidade da AWS](#). Para saber mais sobre os programas de conformidade que se aplicam ao Amazon EMR Sem Servidor, consulte [Serviços da AWS no escopo por programa de conformidade](#).
- Segurança na nuvem — Sua responsabilidade é determinada pelo AWS serviço que você usa. Você também é responsável por outros fatores, incluindo a confidencialidade de seus dados, os requisitos da sua empresa e as leis e normas aplicáveis.

Esta documentação ajuda você a entender como aplicar o modelo de responsabilidade compartilhada usando o Amazon EMR Sem Servidor. Os tópicos explicam como configurar o Amazon EMR Serverless e usar outros AWS serviços para atender aos seus objetivos de segurança e conformidade.

Tópicos

- [Práticas recomendadas de segurança do Amazon EMR Sem Servidor](#)
- [Proteção de dados](#)
- [Identity and Access Management \(IAM\) no Amazon EMR Sem Servidor](#)
- [Propagação de identidades confiáveis](#)
- [Usar o Lake Formation com o EMR Sem Servidor](#)
- [Criptografia entre trabalhadores](#)
- [Criptografia de disco com KMS CMK](#)
- [Secrets Manager para proteção de dados com o EMR Sem Servidor](#)
- [Uso da Concessão de Acesso do Amazon S3 com o EMR Sem Servidor](#)

- [Registro de chamadas de API sem servidor do Amazon EMR usando AWS CloudTrail](#)
- [Validação de conformidade do Amazon EMR Sem Servidor](#)
- [Resiliência no Amazon EMR Sem Servidor](#)
- [Segurança da infraestrutura no Amazon EMR Sem Servidor](#)
- [Análise de configuração e vulnerabilidade no Amazon EMR Sem Servidor](#)

Práticas recomendadas de segurança do Amazon EMR Sem Servidor

O Amazon EMR Sem Servidor disponibiliza diversos recursos de segurança a serem considerados no desenvolvimento e na implementação das suas próprias políticas de segurança. As práticas recomendadas a seguir são diretrizes gerais e não representam uma solução completa de segurança. Como essas práticas recomendadas podem não ser adequadas ou suficientes para o seu ambiente, trate-as como considerações úteis em vez de prescrições.

Aplicação do princípio de privilégio mínimo

O Amazon EMR Sem Servidor fornece uma política de acesso granular para aplicações que usam perfis do IAM, como perfis de execução. Sugerimos que os perfis de execução recebam somente o conjunto de privilégios mínimos obrigatórios para o trabalho, como a cobertura da aplicação e o acesso ao destino do log. Também recomendamos auditar regularmente as permissões dos trabalhos e após qualquer alteração no código da aplicação.

Isolamento de código de uma aplicação não confiável

O EMR Sem Servidor cria isolamento total da rede entre trabalhos pertencentes a diferentes aplicações do EMR Sem Servidor. Nos casos em que o isolamento no nível do trabalho é desejado, considere isolar os trabalhos em diferentes aplicações do EMR Sem Servidor.

Permissões de controle de acesso por perfil (RBAC)

Os administradores devem controlar rigorosamente as permissões de controle de acesso por perfil (RBAC) nas aplicações do EMR Sem Servidor.

Proteção de dados

O [modelo de responsabilidade AWS compartilhada](#) se aplica à proteção de dados no Amazon EMR Serverless. Conforme descrito neste modelo, AWS é responsável por proteger a infraestrutura global que executa toda a AWS nuvem. Você é responsável por manter o controle sobre o conteúdo hospedado nessa infraestrutura. Esse conteúdo inclui as tarefas de configuração e gerenciamento de segurança dos AWS serviços que você usa. Para obter mais informações sobre a privacidade de dados, consulte as [Perguntas frequentes sobre privacidade de dados](#). Para obter informações sobre proteção de dados na Europa, consulte [o Modelo de Responsabilidade AWS Compartilhada e a postagem do blog sobre o GDPR](#) no Blog AWS de Segurança.

Para fins de proteção de dados, sugerimos que você proteja as credenciais da AWS conta e configure contas individuais com AWS Identity and Access Management (IAM). Dessa maneira, cada usuário receberá apenas as permissões necessárias para cumprir suas obrigações de trabalho. Recomendamos também que você proteja seus dados das seguintes formas:

- Use uma autenticação multifator (MFA) com cada conta.
- Use SSL/TLS para se comunicar com AWS os recursos. Recomendamos o TLS 1.2 ou posterior.
- Configure a API e o registro de atividades do usuário com AWS CloudTrail.
- Use soluções AWS de criptografia, juntamente com todos os controles de segurança padrão nos AWS serviços.
- Use serviços gerenciados de segurança avançada, como o Amazon Macie, que ajuda a localizar e proteger dados pessoais armazenados no Amazon S3.
- Use as opções de criptografia do Amazon EMR Sem Servidor para criptografar dados em repouso e em trânsito.
- Se você precisar de módulos criptográficos validados pelo FIPS 140-2 ao acessar AWS por meio de uma interface de linha de comando ou de uma API, use um endpoint FIPS. Para ter mais informações sobre endpoints do FIPS disponíveis, consulte [Federal Information Processing Standard \(FIPS\) 140-2](#).

É altamente recomendável que você nunca coloque informações de identificação confidenciais, como números de conta dos seus clientes, em campos de formato livre, como um campo Nome. Isso inclui quando você trabalha com o Amazon EMR Serverless ou outros AWS serviços usando o console, a API ou. AWS CLI AWS SDKs Todos os dados que você insere no Amazon EMR Sem Servidor ou em outros serviços podem ser separados para inclusão em logs de diagnóstico. Ao

fornecer um URL para um servidor externo, não inclua informações de credenciais no URL para validar a solicitação a esse servidor.

Criptografia em repouso

A criptografia de dados ajuda a impedir que usuários não autorizados leiam dados em um cluster e em sistemas de armazenamento físico de dados associados. Isso inclui dados salvos em mídias persistentes, conhecidos como dados em repouso, e dados que podem ser interceptados enquanto viajam pela rede, conhecidos como dados em trânsito.

A criptografia de dados requer chaves e certificados. Você pode escolher entre várias opções, incluindo chaves gerenciadas por AWS Key Management Service, chaves gerenciadas pelo Amazon S3 e chaves e certificados de fornecedores personalizados fornecidos por você. Ao usar AWS KMS como seu provedor de chaves, cobranças se aplicam pelo armazenamento e uso de chaves de criptografia. Para obter mais informações, consulte [Preços do AWS KMS](#).

Antes de especificar as opções de criptografia, decida quais sistemas de gerenciamento de chaves e de certificados você deseja usar. Em seguida, crie as chaves e os certificados para os provedores personalizados especificados como parte das configurações de criptografia.

Criptografia em repouso para dados do EMRFS no Amazon S3

Cada aplicação do EMR Sem Servidor usa uma versão de lançamento específica, que inclui o EMRFS (EMR File System). A criptografia do Amazon S3 funciona com objetos do Sistema de Arquivos do EMR (EMRFS) lidos e gravados no Amazon S3. Você pode especificar a criptografia do lado do servidor (SSE) ou a criptografia do lado do cliente (CSE) do Amazon S3 como o modo de criptografia padrão ao habilitar a criptografia em repouso. Opcionalmente, especifique diferentes métodos de criptografia para buckets individuais usando Substituições de criptografia por bucket. Independentemente de a criptografia do Amazon S3 estar habilitada, o Transport Layer Security (TLS) criptografa os objetos do EMRFS em trânsito entre os nós do cluster do EMR e o Amazon S3. Se você usa o CSE do Amazon S3 com chaves gerenciadas pelo cliente, o perfil de execução usado para executar trabalhos em uma aplicação do EMR Sem Servidor deve ter acesso à chave. Para obter informações detalhadas sobre a criptografia do Amazon S3, consulte [Proteger dados com criptografia](#) no Guia do desenvolvedor do Amazon Simple Storage Service.

Note

Quando você usa AWS KMS, cobranças são cobradas pelo armazenamento e uso de chaves de criptografia. Para obter mais informações, consulte [Preços do AWS KMS](#).

Criptografia do lado do servidor do Amazon S3

Todos os buckets do Amazon S3 têm criptografia configurada por padrão, e todos os novos objetos que são carregados em um bucket do S3 são automaticamente criptografados em repouso.

O Amazon S3 criptografa os dados no nível do objeto à medida que os grava no disco e os descriptografa quando eles são acessados. Para obter mais informações sobre a SSE, consulte [Proteger os dados usando criptografia do lado do servidor](#) no Guia do desenvolvedor do Amazon Simple Storage Service.

Você pode escolher entre dois sistemas de gerenciamento de chaves diferentes ao especificar a SSE no Amazon EMR Sem Servidor:

- SSE-S3: o Amazon S3 gerencia as chaves para você. Nenhuma configuração adicional é necessária no EMR Sem Servidor.
- SSE-KMS - Você usa um AWS KMS key para configurar políticas adequadas para o EMR Serverless. Nenhuma configuração adicional é necessária no EMR Sem Servidor.

Para usar a AWS KMS criptografia para dados que você grava no Amazon S3, você tem duas opções ao usar a `StartJobRun` API. Você pode habilitar a criptografia para tudo o que é gravado no Amazon S3 ou a criptografia para dados gravados em um bucket específico. Para obter mais informações sobre a API `StartJobRun`, consulte a [Referência da API do EMR Sem Servidor](#).

Para ativar a AWS KMS criptografia para todos os dados que você grava no Amazon S3, use os seguintes comandos ao chamar a `StartJobRun` API.

```
--conf spark.hadoop.fs.s3.enableServerSideEncryption=true  
--conf spark.hadoop.fs.s3.serverSideEncryption.kms.keyId=<kms_id>
```

Para ativar a AWS KMS criptografia de dados que você grava em um bucket específico, use os comandos a seguir ao chamar a `StartJobRun` API.

```
--conf spark.hadoop.fs.s3.bucket.<amzn-s3-demo-bucket1>.enableServerSideEncryption=true  
--conf spark.hadoop.fs.s3.bucket.<amzn-s3-demo-bucket1>.serverSideEncryption.kms.keyId=<kms-id>
```

A SSE com chaves fornecidas pelo cliente (SSE-C) não está disponível para uso com o EMR Sem Servidor.

Criptografia do lado do cliente do Amazon S3

Com a criptografia do lado do cliente do Amazon S3, a criptografia e a descriptografia do Amazon S3 ocorrem no cliente do EMRFS disponível em todas as versões do Amazon EMR. Os objetos são criptografados antes de serem carregados no Amazon S3 e descriptografados após serem baixados. O provedor especificado por você fornece a chave de criptografia que o cliente usa. O cliente pode usar chaves fornecidas pelo AWS KMS (CSE-KMS) ou uma classe Java personalizada que fornece a chave raiz do lado do cliente (CSE-C). As especificações de criptografia são ligeiramente diferentes entre a CSE-KMS e a CSE-C, dependendo do provedor especificado e dos metadados do objeto que está sendo descriptografado ou criptografado. Se você usa o CSE do Amazon S3 com chaves gerenciadas pelo cliente, o perfil de execução usado para executar trabalhos em uma aplicação do EMR Sem Servidor deve ter acesso à chave. Cobranças adicionais do KMS podem ser aplicadas. Para obter mais informações sobre essas diferenças, consulte [Proteger dados usando a criptografia do lado do cliente](#) no Guia do desenvolvedor do Amazon Simple Storage Service.

Criptografia de disco local

Os dados armazenados em armazenamento temporário são criptografados com chaves de propriedade do serviço usando o algoritmo criptográfico AES-256 padrão do setor.

Gerenciamento de chaves

Você pode configurar o KMS para alternar automaticamente suas chaves do KMS. Isso alterna suas chaves uma vez por ano, enquanto salva as chaves antigas indefinidamente para que seus dados ainda possam ser descriptografados. Para obter mais informações, consulte [Chaves mestras rotativas do cliente](#).

Criptografia em trânsito

Os seguintes recursos de criptografia específicos da aplicação estão disponíveis com o Amazon EMR Sem Servidor:

- Spark
 - Por padrão, a comunicação entre drivers e executores do Spark é autenticada e interna. A comunicação de RPC entre drivers e executores é criptografada.
- Hive
 - A comunicação entre o metastore AWS Glue e os aplicativos EMR Serverless ocorre via TLS.

Você deve permitir somente conexões criptografadas via HTTPS (TLS) usando [a SecureTransport condição aws: nas políticas](#) do Amazon S3 bucket IAM.

Identity and Access Management (IAM) no Amazon EMR Sem Servidor

AWS Identity and Access Management (IAM) é uma ferramenta AWS service (Serviço da AWS) que ajuda o administrador a controlar com segurança o acesso aos AWS recursos. Os administradores do IAM controlam quem pode ser autenticado (conectado) e autorizado (ter permissões) para utilizar os recursos do Amazon EMR Sem Servidor. O IAM é um AWS service (Serviço da AWS) que você pode usar sem custo adicional.

Tópicos

- [Público](#)
- [Autenticação com identidades](#)
- [Gerenciar o acesso usando políticas](#)
- [Como o EMR Sem Servidor funciona com o IAM](#)
- [Uso de perfis vinculados ao serviço para o EMR Sem Servidor](#)
- [Perfis de runtime do trabalho para o Amazon EMR Sem Servidor](#)
- [Exemplos de políticas de acesso de usuários do EMR Sem Servidor](#)
- [Políticas para controle de acesso baseado em etiquetas](#)
- [Exemplos de políticas baseadas em identidade do EMR Sem Servidor](#)
- [Atualizações sem servidor do Amazon EMR para políticas gerenciadas AWS](#)
- [Solução de problemas de identidade e acesso da Amazon EMR Sem Servidor](#)

Público

A forma como você usa AWS Identity and Access Management (IAM) difere com base na sua função:

- Usuário do serviço: solicite permissões ao seu administrador se você não conseguir acessar os atributos (consulte [Solução de problemas de identidade e acesso da Amazon EMR Sem Servidor](#)).
- Administrador do serviço: determine o acesso do usuário e envie solicitações de permissão (consulte [Identity and Access Management \(IAM\) no Amazon EMR Sem Servidor](#))

- Administrador do IAM: escreva políticas para gerenciar o acesso (consulte [Exemplos de políticas baseadas em identidade para o EMR Sem Servidor](#))

Autenticação com identidades

A autenticação é a forma como você faz login AWS usando suas credenciais de identidade. Você deve estar autenticado como usuário do IAM ou assumindo uma função do IAM. Usuário raiz da conta da AWS

Você pode fazer login como uma identidade federada usando credenciais de uma fonte de identidade como AWS IAM Identity Center (IAM Identity Center), autenticação de login único ou credenciais. Google/Facebook Para ter mais informações sobre como fazer login, consulte [Como fazer login em sua Conta da AWS](#) no Guia do usuário do Início de Sessão da AWS .

Para acesso programático, AWS fornece um SDK e uma CLI para assinar solicitações criptograficamente. Para ter mais informações, consulte [AWS Signature Version 4 para solicitações de API](#) no Guia do usuário do IAM.

Conta da AWS usuário root

Ao criar um Conta da AWS, você começa com uma identidade de login chamada usuário Conta da AWS raiz que tem acesso completo a todos Serviços da AWS os recursos. É altamente recomendável não usar o usuário-raiz em tarefas diárias. Consulte as tarefas que exigem credenciais de usuário-raiz em [Tarefas que exigem credenciais de usuário-raiz](#) no Guia do usuário do IAM.

Identidade federada

Como prática recomendada, exija que os usuários humanos usem a federação com um provedor de identidade para acessar Serviços da AWS usando credenciais temporárias.

Uma identidade federada é um usuário do seu diretório corporativo, provedor de identidade da web ou Directory Service que acessa Serviços da AWS usando credenciais de uma fonte de identidade. As identidades federadas assumem funções que oferecem credenciais temporárias.

Para o gerenciamento de acesso centralizado, recomendamos AWS IAM Identity Center. Para saber mais, consulte [O que é o IAM Identity Center?](#) no Guia do usuário do AWS IAM Identity Center .

Usuários e grupos do IAM

Um [usuário do IAM](#) é uma identidade com permissões específicas para uma única pessoa ou aplicação. É recomendável usar credenciais temporárias, em vez de usuários do IAM com credenciais de longo prazo. Para obter mais informações, consulte [Exigir que usuários humanos usem a federação com um provedor de identidade para acessar AWS usando credenciais temporárias](#) no Guia do usuário do IAM.

Um [grupo do IAM](#) especifica um conjunto de usuários do IAM e facilita o gerenciamento de permissões para grandes conjuntos de usuários. Para ter mais informações, consulte [Casos de uso de usuários do IAM](#) no Guia do usuário do IAM.

Perfis do IAM

Uma [perfil do IAM](#) é uma identidade com permissões específicas que oferece credenciais temporárias. Você pode assumir uma função [mudando de um usuário para uma função do IAM \(console\)](#) ou chamando uma operação de AWS API AWS CLI ou. Para saber mais, consulte [Métodos para assumir um perfil](#) no Manual do usuário do IAM.

Os perfis do IAM são úteis para acesso de usuário federado, permissões de usuário do IAM temporárias, acesso entre contas, acesso entre serviços e aplicações em execução no Amazon EC2. Consulte mais informações em [Acesso a recursos entre contas no IAM](#) no Guia do usuário do IAM.

Gerenciar o acesso usando políticas

Você controla o acesso AWS criando políticas e anexando-as a AWS identidades ou recursos. Uma política define permissões quando associada a uma identidade ou recurso. AWS avalia essas políticas quando um diretor faz uma solicitação. A maioria das políticas é armazenada AWS como documentos JSON. Para ter mais informações sobre documentos de política JSON, consulte [Visão geral das políticas JSON](#) no Guia do usuário do IAM.

Por meio de políticas, os administradores especificam quem tem acesso a que, definindo qual entidade principal pode realizar ações em quais recursos e sob quais condições.

Por padrão, usuários e perfis não têm permissões. Um administrador do IAM cria políticas do IAM e as adiciona aos perfis, os quais os usuários podem então assumir. As políticas do IAM definem permissões, independentemente do método usado para realizar a operação.

Políticas baseadas em identidade

As políticas baseadas em identidade são documentos de políticas de permissão JSON que você anexa a uma identidade (usuário, grupo ou perfil). Essas políticas controlam quais ações as identidades podem realizar, em quais recursos e sob quais condições. Para saber como criar uma política baseada em identidade, consulte [Definir permissões personalizadas do IAM com as políticas gerenciadas pelo cliente](#) no Guia do Usuário do IAM.

As políticas baseadas em identidade podem ser políticas em linha (incorporadas diretamente em uma única identidade) ou políticas gerenciadas (políticas autônomas anexadas a várias identidades). Para saber como escolher entre uma política gerenciada e políticas em linha, consulte [Escolher entre políticas gerenciadas e políticas em linha](#) no Guia do usuário do IAM.

Políticas baseadas em recursos

Políticas baseadas em recursos são documentos de políticas JSON que você anexa a um recurso. Entre os exemplos estão políticas de confiança de perfil do IAM e políticas de bucket do Amazon S3. Em serviços compatíveis com políticas baseadas em recursos, os administradores de serviço podem usá-las para controlar o acesso a um recurso específico. É necessário [especificar uma entidade principal](#) em uma política baseada em recursos.

Políticas baseadas em recursos são políticas em linha localizadas nesse serviço. Você não pode usar políticas AWS gerenciadas do IAM em uma política baseada em recursos.

Outros tipos de política

AWS oferece suporte a tipos de políticas adicionais que podem definir o máximo de permissões concedidas por tipos de políticas mais comuns:

- Limites de permissões: definem o número máximo de permissões que uma política baseada em identidade pode conceder a uma entidade do IAM. Para saber mais sobre limites de permissões, consulte [Limites de permissões para identidades do IAM](#) no Guia do usuário do IAM.
- Políticas de controle de serviço (SCPs) — Especifique as permissões máximas para uma organização ou unidade organizacional em AWS Organizations. Para saber mais, consulte [Políticas de controle de serviço](#) no Guia do usuário do AWS Organizations .
- Políticas de controle de recursos (RCPs) — Defina o máximo de permissões disponíveis para recursos em suas contas. Para obter mais informações, consulte [Políticas de controle de recursos \(RCPs\)](#) no Guia AWS Organizations do usuário.

- Políticas de sessão: políticas avançadas transmitidas como um parâmetro durante a criação de uma sessão temporária para um perfil ou um usuário federado. Para saber mais, consulte [Políticas de sessão](#) no Guia do usuário do IAM.

Vários tipos de política

Quando vários tipos de política são aplicáveis a uma solicitação, é mais complicado compreender as permissões resultantes. Para saber como AWS determinar se uma solicitação deve ser permitida quando vários tipos de políticas estão envolvidos, consulte [Lógica de avaliação de políticas](#) no Guia do usuário do IAM.

Como o EMR Sem Servidor funciona com o IAM

Antes de usar o IAM para gerenciar o acesso ao Amazon EMR Sem Servidor, entenda que recursos do IAM estão disponíveis para uso com o Amazon EMR Sem Servidor.

Recursos do IAM que você pode usar com o EMR Sem Servidor

Recurso do IAM	Suporte do Amazon EMR Sem Servidor
Políticas baseadas em identidade	Sim
Políticas baseadas em recurso	Não
Ações de políticas	Sim
Recursos de políticas	Sim
Chaves de condição de políticas	Não
ACLs	Não
ABAC (tags em políticas)	Sim
Credenciais temporárias	Sim
Permissões de entidade principal	Sim
Perfis de serviço	Não
Perfis vinculados ao serviço	Sim

Para obter uma visão de alto nível de como o EMR Serverless e AWS outros serviços funcionam com a maioria dos recursos do IAM, consulte os [serviços que funcionam com AWS o IAM](#) no Guia do usuário do IAM.

Políticas baseadas em identidade do EMR Sem Servidor

Compatível com políticas baseadas em identidade: sim

As políticas baseadas em identidade são documentos de políticas de permissões JSON que podem ser anexados a uma identidade, como usuário do IAM, grupo de usuários ou perfil. Essas políticas controlam quais ações os usuários e perfis podem realizar, em quais recursos e em que condições. Para saber como criar uma política baseada em identidade, consulte [Definir permissões personalizadas do IAM com as políticas gerenciadas pelo cliente](#) no Guia do Usuário do IAM.

Com as políticas baseadas em identidade do IAM, é possível especificar ações e recursos permitidos ou negados, assim como as condições sob as quais as ações são permitidas ou negadas. Para saber mais sobre todos os elementos que podem ser usados em uma política JSON, consulte [Referência de elemento de política JSON do IAM](#) no Guia do usuário do IAM.

Exemplos de políticas baseadas em identidade para o EMR Sem Servidor

Para acessar exemplos de políticas baseadas em identidade do Amazon EMR Sem Servidor, consulte [Exemplos de políticas baseadas em identidade do EMR Sem Servidor](#).

Políticas baseadas em recursos no EMR Sem Servidor

Compatibilidade com políticas baseadas em recursos: não

Políticas baseadas em recursos são documentos de políticas JSON que você anexa a um recurso. São exemplos de políticas baseadas em recursos as políticas de confiança de perfil do IAM e as políticas de bucket do Amazon S3. Em serviços compatíveis com políticas baseadas em recursos, os administradores de serviço podem usá-las para controlar o acesso a um recurso específico. Para o atributo ao qual a política está anexada, a política define quais ações uma entidade principal especificado pode executar nesse atributo e em que condições. É necessário [especificar uma entidade principal](#) em uma política baseada em recursos. Os diretores podem incluir contas, usuários, funções, usuários federados ou. Serviços da AWS

Para permitir o acesso entre contas, é possível especificar uma conta inteira ou as entidades do IAM em outra conta como a entidade principal em uma política baseada em recursos. Consulte mais informações em [Acesso a recursos entre contas no IAM](#) no Guia do usuário do IAM.

Ações de políticas do EMR Sem Servidor

Compatível com ações de políticas: sim

Os administradores podem usar políticas AWS JSON para especificar quem tem acesso ao quê. Ou seja, qual entidade principal pode executar ações em quais recursos e em que condições.

O elemento `Action` de uma política JSON descreve as ações que podem ser usadas para permitir ou negar acesso em uma política. Incluem ações em uma política para conceder permissões para executar a operação associada.

Para ver uma lista de ações do EMR Sem Servidor, consulte [Ações, recursos e chaves de condição para o Amazon EMR Sem Servidor](#) na Referência de autorização do serviço.

As ações de políticas no EMR Sem Servidor usam o prefixo a seguir antes da ação.

```
emr-serverless
```

Para especificar várias ações em uma única declaração, separe-as com vírgulas.

```
"Action": [  
  "emr-serverless:action1",  
  "emr-serverless:action2"  
]
```

Para acessar exemplos de políticas baseadas em identidade do Amazon EMR Sem Servidor, consulte [Exemplos de políticas baseadas em identidade do EMR Sem Servidor](#).

Recursos de políticas do EMR Sem Servidor

Compatível com recursos de políticas: sim

Os administradores podem usar políticas AWS JSON para especificar quem tem acesso ao quê. Ou seja, qual entidade principal pode executar ações em quais recursos e em que condições.

O elemento de política JSON `Resource` especifica o objeto ou os objetos aos quais a ação se aplica. Como prática recomendada, especifique um recurso usando seu [nome do recurso da Amazon](#)

[\(ARN\)](#). Para ações que não oferecem compatibilidade com permissões em nível de recurso, use um curinga (*) para indicar que a instrução se aplica a todos os recursos.

```
"Resource": "*"
```

Para consultar uma lista dos tipos de recursos sem servidor do Amazon EMR e seus ARNs, consulte Recursos [definidos pelo Amazon EMR Serverless](#) na Referência de autorização de serviço. Para saber com quais ações especificar o ARN de cada recurso, consulte [Ações, recursos e chaves de condição para Amazon EMR Sem Servidor](#).

Para acessar exemplos de políticas baseadas em identidade do Amazon EMR Sem Servidor, consulte [Exemplos de políticas baseadas em identidade do EMR Sem Servidor](#).

Chaves de condição de política do EMR Sem Servidor

Suporte para chaves de condições de política

Suporta chaves de condição de política específicas de serviço	Não
---------------------------------------------------------------	-----

Os administradores podem usar políticas AWS JSON para especificar quem tem acesso ao quê. Ou seja, qual entidade principal pode executar ações em quais recursos e em que condições.

O elemento `Condition` especifica quando as instruções são executadas com base em critérios definidos. É possível criar expressões condicionais que usem [agentes de condição](#), como “igual a” ou “menor que”, para fazer a condição da política corresponder aos valores na solicitação. Para ver todas as chaves de condição AWS globais, consulte as [chaves de contexto de condição AWS global](#) no Guia do usuário do IAM.

Para obter uma lista de chaves de condição do Amazon EMR Sem Servidor, das ações e dos recursos que você pode usar, consulte [Actions, resources, and condition keys for Amazon EMR Sem Servidor](#) na Referência de autorização do serviço.

Todas as ações do Amazon EC2 oferecem suporte às chaves de condição `aws:RequestedRegion` e `ec2:Region`. Para obter mais informações, consulte [Exemplo: Restringir o acesso a uma região específica](#).

Listas de controle de acesso (ACLs) no EMR Serverless

Suportes ACLs: Não

As listas de controle de acesso (ACLs) controlam quais diretores (membros da conta, usuários ou funções) têm permissões para acessar um recurso. ACLs são semelhantes às políticas baseadas em recursos, embora não usem o formato de documento de política JSON.

Controle de acesso por atributo (ABAC) com o EMR Sem Servidor

Suporte ao controle de acesso por atributo (ABAC)

Oferece compatibilidade com ABAC (tags em políticas)	Sim
------------------------------------------------------	-----

O controle de acesso por atributo (ABAC) é uma estratégia de autorização que define permissões com base em atributos chamados de tags. Você pode anexar tags a entidades e AWS recursos do IAM e, em seguida, criar políticas ABAC para permitir operações quando a tag do diretor corresponder à tag no recurso.

Para controlar o acesso baseado em tags, forneça informações sobre as tags no [elemento de condição](#) de uma política usando as `aws:ResourceTag/key-name`, `aws:RequestTag/key-name` ou chaves de condição `aws:TagKeys`.

Se um serviço for compatível com as três chaves de condição para cada tipo de recurso, o valor será Sim para o serviço. Se um serviço for compatível com as três chaves de condição somente para alguns tipos de recursos, o valor será Parcial

Para saber mais sobre o ABAC, consulte [Definir permissões com autorização do ABAC](#) no Guia do usuário do IAM. Para visualizar um tutorial com etapas para configurar o ABAC, consulte [Usar controle de acesso por atributo \(ABAC\)](#) no Guia do usuário do IAM.

Uso de credenciais temporárias com o EMR Sem Servidor

Compatível com credenciais temporárias: sim

As credenciais temporárias fornecem acesso de curto prazo aos AWS recursos e são criadas automaticamente quando você usa a federação ou troca de funções. AWS recomenda que você gere credenciais temporárias dinamicamente em vez de usar chaves de acesso de longo prazo. Para ter

mais informações, consulte [Credenciais de segurança temporárias no IAM](#) e [Serviços da Serviços da AWS que funcionam com o IAM](#) no Guia do usuário do IAM.

Permissões de entidades principais entre serviços do EMR Sem Servidor

Compatibilidade com o recurso de encaminhamento de sessões de acesso (FAS): sim

As sessões de acesso direto (FAS) usam as permissões do principal chamando um AWS service (Serviço da AWS), combinadas com a solicitação AWS service (Serviço da AWS) de fazer solicitações aos serviços posteriores. Para obter detalhes da política ao fazer solicitações de FAS, consulte [Sessões de acesso direto](#).

Perfis de serviço do EMR Sem Servidor

Oferece suporte a perfis de serviço	Não
-------------------------------------	-----

Perfis vinculados ao serviço do EMR Sem Servidor

Oferece suporte a perfis vinculados ao serviço	Sim
------------------------------------------------	-----

Para obter detalhes sobre como criar ou gerenciar perfis vinculados a serviços, consulte [Serviços da AWS que funcionam com o IAM](#). Encontre um serviço na tabela que inclua um Yes na coluna Perfil vinculado ao serviço. Escolha o link Sim para acessar a documentação do perfil vinculado a serviço desse serviço.

Uso de perfis vinculados ao serviço para o EMR Sem Servidor

[O Amazon EMR Serverless usa funções vinculadas a serviços AWS Identity and Access Management \(IAM\)](#). Um perfil vinculado ao serviço é um tipo exclusivo de perfil do IAM vinculado diretamente ao EMR Sem Servidor. As funções vinculadas ao serviço são predefinidas pelo EMR Serverless e incluem todas as permissões que o serviço exige para chamar outros serviços em seu nome. AWS

Um perfil vinculado ao serviço facilita a configuração do EMR Sem Servidor porque você não precisa adicionar as permissões necessárias manualmente. O EMR Sem Servidor define as permissões desses perfis vinculados ao serviço e, exceto se definido de outra forma, somente o EMR Sem

Servidor pode assumir os próprios perfis. As permissões definidas incluem a política de confiança e a política de permissões, que não pode ser anexada a nenhuma outra entidade do IAM.

Um perfil vinculado ao serviço poderá ser excluído somente após excluir seus atributos relacionados. Isso protege seus recursos do EMR Sem Servidor, porque você não pode remover a permissão por engano para acessá-los.

Para obter informações sobre outros serviços compatíveis com perfis vinculados a serviços, consulte [Serviços da AWS que funcionam com o IAM](#) e procure os serviços que apresentam Sim na coluna Perfis vinculados a serviços. Escolha um Sim com um link para acessar a documentação do perfil vinculado para esse serviço.

Permissões de perfil vinculado ao serviço do EMR Sem Servidor

O EMR Serverless usa a função vinculada ao serviço chamada `AWSServiceRoleForAmazonEMRServerless` para permitir que ele ligue em seu nome. AWS APIs

A função `AWSServiceRoleForAmazonEMRServerless` vinculada ao serviço confia nos seguintes serviços para assumir a função:

- `ops.emr-serverless.amazonaws.com`

A política de permissões do perfil chamada `AmazonEMRServerlessServiceRolePolicy` permite que o EMR Sem Servidor conclua as ações a seguir nos recursos especificados.

Note

Como o conteúdo da política gerenciada muda, a política mostrada aqui pode estar desatualizada. Veja a maioria das up-to-date políticas [da Amazon EMRServerless ServiceRolePolicy](#) no Console de gerenciamento da AWS.

- Ação: `ec2:CreateNetworkInterface`
- Ação: `ec2>DeleteNetworkInterface`
- Ação: `ec2:DescribeNetworkInterfaces`
- Ação: `ec2:DescribeSecurityGroups`
- Ação: `ec2:DescribeSubnets`
- Ação: `ec2:DescribeVpcs`

- Ação: ec2:DescribeDhcpOptions
- Ação: ec2:DescribeRouteTables
- Ação: cloudwatch:PutMetricData

A política a seguir é a AmazonEMRServerlessServiceRolePolicy completa.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "EC2PolicyStatement",
      "Effect": "Allow",
      "Action": [
        "ec2:CreateNetworkInterface",
        "ec2>DeleteNetworkInterface",
        "ec2:DescribeNetworkInterfaces",
        "ec2:DescribeSecurityGroups",
        "ec2:DescribeSubnets",
        "ec2:DescribeVpcs",
        "ec2:DescribeDhcpOptions",
        "ec2:DescribeRouteTables"
      ],
      "Resource": [
        "*"
      ]
    },
    {
      "Sid": "CloudWatchPolicyStatement",
      "Effect": "Allow",
      "Action": [
        "cloudwatch:PutMetricData"
      ],
      "Resource": [
        "*"
      ],
      "Condition": {
        "StringEquals": {
          "cloudwatch:namespace": [
            "AWS/EMRServerless",

```

```

    "AWS/Usage"
  ]
}
]
}

```

A política de confiança a seguir está anexada a esse perfil para permitir que a entidade principal do EMR Sem Servidor assuma o perfil.

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "sts:AssumeRole"
      ],
      "Resource": "arn:aws:iam::123456789012:role/aws-service-role/emr-serverless.amazonaws.com/AWSServiceRoleForEMRServerless",
      "Sid": "AllowSTSAssumerole"
    }
  ]
}

```

É necessário configurar as permissões para permitir que uma entidade do IAM (como um usuário, grupo ou perfil) crie, edite ou exclua uma função vinculada ao serviço. Para obter mais informações, consulte [Permissões de perfil vinculado a serviços](#) no Guia do usuário do IAM.

Criação de um perfil vinculado ao serviço do EMR Sem Servidor

Não é necessário criar manualmente um perfil vinculado ao serviço. Quando você cria um novo aplicativo EMR Serverless no (Console de gerenciamento da AWS usando o EMR Studio), no AWS CLI ou na API, o EMR Serverless cria a função AWS vinculada ao serviço para você. É necessário configurar as permissões para permitir que uma entidade do IAM (como um usuário, grupo ou perfil) crie, edite ou exclua uma função vinculada ao serviço.

Para criar a função `AWSServiceRoleForAmazonEMRServerless` vinculada ao serviço usando o IAM

Adicione a instrução a seguir à política de permissões da entidade do IAM que precisa criar o perfil vinculado ao serviço.

```
{
  "Effect": "Allow",
  "Action": [
    "iam:CreateServiceLinkedRole"
  ],
  "Resource": "arn:aws:iam::*:role/aws-service-role/ops.emr-serverless.amazonaws.com/AWSServiceRoleForAmazonEMRServerless*",
  "Condition": {"StringLike": {"iam:AWSServiceName": "ops.emr-serverless.amazonaws.com"}}
}
```

Se você excluir esse perfil vinculado ao serviço e precisar criá-lo novamente, use esse mesmo processo para recriar o perfil na sua conta. Ao criar uma aplicação do EMR Sem Servidor, o EMR Sem Servidor cria o perfil vinculado ao serviço para você novamente.

Você também pode usar o console do IAM para criar um perfil vinculado ao serviço com o caso de uso EMR Sem Servidor. Na AWS CLI ou na AWS API, crie uma função vinculada ao serviço com o nome do `ops.emr-serverless.amazonaws.com` serviço. Para obter mais informações, consulte [Criar um perfil vinculado ao serviço](#) no Guia do usuário do IAM. Se você excluir essa função vinculada ao serviço, use esse mesmo processo para criar a função novamente.

Edição de um perfil vinculado ao serviço do EMR Sem Servidor

O EMR Serverless não permite que você edite a função `AWSServiceRoleForAmazonEMRServerless` vinculada ao serviço porque várias entidades podem fazer referência à função. Você não pode editar a política do IAM AWS de propriedade que a função vinculada ao serviço do EMR Serverless usa, pois ela contém todas as permissões necessárias que o EMR Serverless precisa. No entanto, você poderá editar a descrição do perfil usando o IAM.

Para editar a descrição da função `AWSServiceRoleForAmazonEMRServerless` vinculada ao serviço usando o IAM

Adicione a instrução abaixo à política de permissões da entidade do IAM para a qual precise editar a descrição de um perfil vinculado ao serviço.

```
{
```

```

    "Effect": "Allow",
    "Action": [
        "iam: UpdateRoleDescription"
    ],
    "Resource": "arn:aws:iam::*:role/aws-service-role/ops.emr-serverless.amazonaws.com/
AWSServiceRoleForAmazonEMRServerless*",
    "Condition": {"StringLike": {"iam:AWSServiceName": "ops.emr-
serverless.amazonaws.com"}}
}

```

Para obter mais informações, consulte [Editar um perfil vinculado ao serviço](#) no Guia do usuário do IAM.

Exclusão de um perfil vinculado a serviço do EMR Sem Servidor

Se você não precisar mais usar um recurso ou serviço que requer um perfil vinculado ao serviço, é recomendável excluí-lo. Dessa forma, você não tem uma entidade não utilizada que não seja monitorada ativamente ou mantida. Contudo, você deve excluir todas as aplicações do EMR Sem Servidor em todas as regiões para poder excluir o perfil vinculado ao serviço.

Note

Se o serviço EMR Sem Servidor estiver usando o perfil quando você tenta excluir os recursos associados ao perfil, a exclusão poderá falhar. Se isso acontecer, espere alguns minutos e tente a operação novamente.

Para excluir a função `AWSService RoleForAmazon EMRServerless` vinculada ao serviço usando o IAM

Adicione a instrução a seguir à política de permissões da entidade do IAM que precisa excluir um perfil vinculado ao serviço.

```

{
    "Effect": "Allow",
    "Action": [
        "iam:DeleteServiceLinkedRole",
        "iam:GetServiceLinkedRoleDeletionStatus"
    ],
    "Resource": "arn:aws:iam::*:role/aws-service-role/ops.emr-serverless.amazonaws.com/
AWSServiceRoleForAmazonEMRServerless*",
}

```

```
"Condition": {"StringLike": {"iam:AWSServiceName": "ops.emr-serverless.amazonaws.com"}}
}
```

Como excluir manualmente o perfil vinculado ao serviço usando o IAM

Use o console do IAM AWS CLI, o ou a AWS API para excluir a função AWSService RoleForAmazon EMRServerless vinculada ao serviço. Para obter mais informações, consulte [Excluir um perfil vinculado ao serviço](#) no Guia do usuário do IAM.

Regiões compatíveis com perfis vinculados ao serviço do EMR Sem Servidor

O EMR Sem Servidor é compatível com a utilização de perfis vinculados ao serviço em todas as regiões em que o serviço está disponível. Para obter mais informações, consulte [Regiões e endpoints da AWS](#).

Perfis de runtime do trabalho para o Amazon EMR Sem Servidor

Você pode especificar as permissões do perfil do IAM que a execução de trabalho do EMR Sem Servidor pode assumir ao chamar outros serviços em seu nome. Isso inclui acesso ao Amazon S3 para quaisquer fontes de dados, destinos e outros AWS recursos, como clusters do Amazon Redshift e tabelas do DynamoDB. Para saber mais sobre como criar um perfil, consulte [Criação de um perfil de runtime de trabalhos](#).

Exemplos de políticas de runtime

Você pode anexar uma política de runtime, como a mostrada a seguir, a um perfil de runtime. A seguinte política de runtime de trabalhos permite:

- Acesso de leitura aos buckets do Amazon S3 com exemplos do EMR.
- Acesso total aos buckets do S3.
- Crie e leia o acesso ao AWS Glue Data Catalog.

Para adicionar acesso a outros AWS recursos, como o DynamoDB, você precisará incluir permissões para eles na política ao criar a função de tempo de execução.

JSON

```
{
```

```
"Version": "2012-10-17",
"Statement": [
  {
    "Sid": "ReadAccessForEMRSamples",
    "Effect": "Allow",
    "Action": [
      "s3:GetObject",
      "s3:ListBucket"
    ],
    "Resource": [
      "arn:aws:s3::*.elasticmapreduce",
      "arn:aws:s3::*.elasticmapreduce/*"
    ]
  },
  {
    "Sid": "FullAccessToS3Bucket",
    "Effect": "Allow",
    "Action": [
      "s3:PutObject",
      "s3:GetObject",
      "s3:ListBucket",
      "s3:DeleteObject"
    ],
    "Resource": [
      "arn:aws:s3:::amzn-s3-demo-bucket",
      "arn:aws:s3:::amzn-s3-demo-bucket/*"
    ]
  },
  {
    "Sid": "GlueCreateAndReadDataCatalog",
    "Effect": "Allow",
    "Action": [
      "glue:GetDatabase",
      "glue:CreateDatabase",
      "glue:GetDataBases",
      "glue:CreateTable",
      "glue:GetTable",
      "glue:UpdateTable",
      "glue>DeleteTable",
      "glue:GetTables",
      "glue:GetPartition",
      "glue:GetPartitions",
      "glue:CreatePartition",
      "glue:BatchCreatePartition",
```

```

    "glue:GetUserDefinedFunctions"
  ],
  "Resource": [
    "*"
  ]
}
]
}

```

Transmissão de privilégios de perfil

Você pode anexar políticas de permissões do IAM ao perfil de um usuário para permitir que ele passe apenas perfis aprovados. Isso permite que os administradores controlem quais usuários podem transmitir perfis específicos de runtime de trabalho para trabalhos do EMR Sem Servidor. Para saber mais sobre como definir permissões, consulte [Conceder permissões a um usuário para passar uma função para um AWS serviço](#).

Confira a seguir um exemplo de política que permite transmitir um perfil de runtime de trabalhos para a entidade principal de serviço do EMR Sem Servidor.

```

{
  "Effect": "Allow",
  "Action": "iam:PassRole",
  "Resource": "arn:aws:iam::1234567890:role/JobRuntimeRoleForEMRServerless",
  "Condition": {
    "StringLike": {
      "iam:PassedToService": "emr-serverless.amazonaws.com"
    }
  }
}

```

Políticas de permissão gerenciada associadas a perfis de runtime

Ao enviar execuções de trabalhos para o EMR Sem Servidor por meio do console do EMR Studio, há uma etapa em que você escolhe um Perfil de runtime para associar à sua aplicação. Há políticas gerenciadas subjacentes associadas a cada seleção no console que é importante conhecer. As três seleções são:

1. Todos os buckets — Quando você escolhe essa opção, ela especifica a política FullAccess AWS gerenciada do [AmazonS3](#), que fornece acesso total a todos os buckets.

2. Buckets específicos: especifica o identificador de Nome de recurso da Amazon (ARN) de cada bucket que você escolhe. Não há uma política gerenciada subjacente incluída.
3. Nenhuma: nenhuma permissão de política gerenciada está incluída.

Sugerimos adicionar buckets específicos. Se você escolher todos os buckets, lembre-se de que isso define o acesso total a todos os buckets.

Exemplos de políticas de acesso de usuários do EMR Sem Servidor

Você pode configurar políticas refinadas para seus usuários, dependendo das ações que cada um deve executar ao interagir com aplicações do EMR Sem Servidor. As políticas a seguir são exemplos que podem ajudar na configuração das permissões apropriadas para os usuários. Esta seção se concentra somente nas políticas do EMR Sem Servidor. Para obter exemplos de políticas de usuário do EMR Studio, consulte [Configurar permissões de usuário do EMR Studio](#). Para obter informações sobre como anexar políticas aos usuários do IAM (entidades principais), consulte [Gerenciar políticas do IAM](#) no Guia do usuário do IAM.

Política de usuários avançados

Para conceder todas as ações necessárias ao EMR Sem Servidor, crie e anexe uma política AmazonEMRServerlessFullAccess ao usuário, perfil ou grupo do IAM necessário.

Confira a seguir um exemplo de política que permite que usuários avançados criem e modifiquem aplicações do EMR Sem Servidor, além de realizar outras ações, como enviar e depurar trabalhos. Ele revela todas as ações que o EMR Sem Servidor requer para outros serviços.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "EMRServerlessActions",
      "Effect": "Allow",
      "Action": [
        "emr-serverless:CreateApplication",
        "emr-serverless:UpdateApplication",
        "emr-serverless>DeleteApplication",
        "emr-serverless:ListApplications",
```

```

    "emr-serverless:GetApplication",
    "emr-serverless:StartApplication",
    "emr-serverless:StopApplication",
    "emr-serverless:StartJobRun",
    "emr-serverless:CancelJobRun",
    "emr-serverless:ListJobRuns",
    "emr-serverless:GetJobRun"
  ],
  "Resource": [
    "*"
  ]
}
]
}

```

Ao habilitar a conectividade de rede com a VPC, as aplicações do EMR Sem Servidor criam interfaces de rede elástica (ENIs) do Amazon EC2 para se comunicar com os recursos da VPC. A política a seguir garante que qualquer novo EC2 ENIs seja criado somente no contexto dos aplicativos EMR Serverless.

Note

É altamente recomendável definir essa política para garantir que os usuários não possam criar o EC2, ENIs exceto no contexto da inicialização de aplicativos EMR Serverless.

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "AllowEC2ENICreationWithEMRTags",
      "Effect": "Allow",
      "Action": [
        "ec2:CreateNetworkInterface"
      ],
      "Resource": [
        "arn:aws:ec2:*:*:network-interface/*"
      ],
    }
  ],
}

```

```

    "Condition": {
      "StringEquals": {
        "aws:CalledViaLast": "ops.emr-serverless.amazonaws.com"
      }
    }
  ]
}

```

Se quiser restringir o acesso do EMR Sem Servidor a determinadas sub-redes, você pode marcar cada sub-rede com uma condição de tag. Essa política do IAM garante que os aplicativos EMR Serverless só possam criar ENIs EC2 dentro de sub-redes permitidas.

```

{
  "Sid": "AllowEC2ENICreationInSubnetAndSecurityGroupWithEMRTags",
  "Effect": "Allow",
  "Action": [
    "ec2:CreateNetworkInterface"
  ],
  "Resource": [
    "arn:aws:ec2:*:*:subnet/*",
    "arn:aws:ec2:*:*:security-group/*"
  ],
  "Condition": {
    "StringEquals": {
      "aws:ResourceTag/KEY": "VALUE"
    }
  }
}

```

Important

Se você for um administrador ou usuário avançado criando sua primeira aplicação, deverá configurar suas políticas de permissão para permitir a criação de um perfil vinculado a serviços do EMR Sem Servidor. Consulte para saber mais: [Uso de perfis vinculados ao serviço para o EMR Sem Servidor](#).

A política do IAM a seguir permite criar um perfil vinculado ao serviço do EMR Sem Servidor para a sua conta.

```
{
  "Sid": "AllowEMRServerlessServiceLinkedRoleCreation",
  "Effect": "Allow",
  "Action": "iam:CreateServiceLinkedRole",
  "Resource": "arn:aws:iam::account-id:role/aws-service-role/ops.emr-serverless.amazonaws.com/AWSServiceRoleForAmazonEMRServerless"
}
```

Política de engenheiro de dados

O exemplo a seguir é de uma política que permite aos usuários permissões somente leitura em aplicações do EMR Sem Servidor, bem como a capacidade de enviar e depurar trabalhos. Lembre-se de que, como essa política não nega explicitamente as ações, uma declaração de política diferente ainda pode ser usada para conceder acesso a ações específicas.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "EMRServerlessActions",
      "Effect": "Allow",
      "Action": [
        "emr-serverless:ListApplications",
        "emr-serverless:GetApplication",
        "emr-serverless:StartApplication",
        "emr-serverless:StartJobRun",
        "emr-serverless:CancelJobRun",
        "emr-serverless:ListJobRuns",
        "emr-serverless:GetJobRun"
      ],
      "Resource": [
        "*"
      ]
    }
  ]
}
```

Uso de tags para controle de acesso com

Você pode usar condições de tag para controle de acesso refinado. Por exemplo, você pode restringir usuários de uma equipe para que eles só possam enviar trabalhos a aplicações do EMR Sem Servidor marcados com o nome da equipe deles.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "EMRServerlessActions",
      "Effect": "Allow",
      "Action": [
        "emr-serverless:ListApplications",
        "emr-serverless:GetApplication",
        "emr-serverless:StartApplication",
        "emr-serverless:StartJobRun",
        "emr-serverless:CancelJobRun",
        "emr-serverless:ListJobRuns",
        "emr-serverless:GetJobRun"
      ],
      "Resource": [
        "*"
      ]
    }
  ]
}
```

Políticas para controle de acesso baseado em etiquetas

Você pode usar condições na política baseada em identidade para controlar o acesso a aplicações e execuções de trabalhos com base em tags.

Os exemplos a seguir demonstram diferentes cenários e maneiras de usar operadores de condição com chaves de condição do EMR Sem Servidor. Estas instruções de política do IAM são destinadas somente para fins de demonstração e não devem ser usadas em ambientes de produção. Há várias maneiras de combinar declarações de políticas para conceder e negar permissões de acordo com

seus requisitos. Para obter mais informações sobre como planejar e testar políticas do IAM, consulte o [Guia do usuário do IAM](#).

Important

Recusar, explicitamente, permissões para ações de uso de tags é uma consideração importante. Isso evita que os usuários façam a marcação de um recurso e, assim, concedam a si mesmos permissões que você não pretendia conceder. Se as ações de marcação de um recurso não forem negadas, um usuário poderá modificar as etiquetas e driblar a intenção das políticas baseadas em etiquetas. Para obter um exemplo de política que nega ações de marcação, consulte [Negação de acesso para adição ou remoção de etiquetas](#).

Os exemplos a seguir demonstram a políticas de permissões baseadas em identidade que são usadas para controlar as ações permitidas com aplicações do EMR Sem Servidor.

Ações permitidas somente em recursos com valores de etiquetas específicos

No exemplo de política a seguir, o operador de condição `StringEquals` tenta corresponder `dev` com o valor da tag `department`. Se a tag "Department" não tiver sido adicionada à aplicação ou não contiver o valor `dev`, a política não se aplicará e as ações não serão permitidas por essa política. Se nenhuma outra instrução de política permitir as ações, o usuário só poderá trabalhar com aplicações que tenham essa tag com esse valor.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "emr-serverless:GetApplication"
      ],
      "Resource": [
        "*"
      ],
      "Condition": {
        "StringEquals": {
          "aws:ResourceTag/department": "dev"
        }
      }
    }
  ]
}
```

```

    }
  },
  "Sid": "AllowEMRSERVERLESSGetapplication"
}
]
}

```

Você também pode especificar vários valores de tag usando um operador de condição. Por exemplo, para permitir ações em aplicações em que a tag `department` contenha o valor `dev` ou `test`, substitua o bloco condicional no exemplo anterior pelo conteúdo a seguir.

```

"Condition": {
  "StringEquals": {
    "emr-serverless:ResourceTag/department": ["dev", "test"]
  }
}

```

Marcação obrigatória na criação de um recurso

No exemplo abaixo, a tag precisa ser aplicada ao criar a aplicação.

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "emr-serverless:CreateApplication"
      ],
      "Resource": [
        "*"
      ],
      "Condition": {
        "StringEquals": {
          "aws:RequestedRegion": "us-east-1"
        }
      }
    },
    {
      "Sid": "AllowEMRSERVERLESSCreateapplication"
    }
  ]
}

```

```

]
}

```

A instrução de política apresentada a seguir permite que um usuário crie uma aplicação somente se ela tiver uma tag `department`, que pode conter qualquer valor.

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "emr-serverless:CreateApplication"
      ],
      "Resource": [
        "*"
      ],
      "Condition": {
        "StringEquals": {
          "aws:RequestedRegion": ["us-east-1", "us-west-2"]
        }
      },
      "Sid": "AllowEMRSERVERLESSCreateapplication"
    }
  ]
}

```

Negação de acesso para adição ou remoção de etiquetas

Essa política impede que um usuário adicione ou remova tags em aplicações do EMR Sem Servidor com uma tag `department` cujo valor não seja `dev`.

JSON

```

{
  "Version": "2012-10-17",
  "Statement": [

```

```
{
  "Effect": "Deny",
  "Action": [
    "emr-serverless:TagResource",
    "emr-serverless:UntagResource"
  ],
  "Resource": [
    "*"
  ],
  "Condition": {
    "StringNotEquals": {
      "aws:PrincipalTag/department": "dev"
    }
  },
  "Sid": "AllowEMRSERVERLESSTagresource"
}
```

Exemplos de políticas baseadas em identidade do EMR Sem Servidor

Por padrão, os usuários e os perfis não têm permissão para criar ou modificar os recursos do Amazon EMR Sem Servidor. Para conceder permissão aos usuários para executar ações nos recursos que eles precisam, um administrador do IAM pode criar políticas do IAM.

Para aprender a criar uma política baseada em identidade do IAM ao usar esses documentos de política em JSON de exemplo, consulte [Criar políticas do IAM \(console\)](#) no Guia do usuário do IAM.

Para obter detalhes sobre ações e tipos de recursos definidos pelo Amazon EMR Serverless, incluindo o formato de cada um dos tipos de recursos, consulte [Ações, recursos e chaves de condição ARNs para o Amazon EMR Serverless](#) na Referência de autorização de serviço.

Tópicos

- [Práticas recomendadas de política](#)
- [Permitir que os usuários acessem suas próprias permissões](#)

Práticas recomendadas de política

Note

O EMR Sem Servidor não oferece suporte a políticas gerenciadas, portanto, a primeira prática listada abaixo não se aplica.

As políticas baseadas em identidade determinam se alguém pode criar, acessar ou excluir recursos do Amazon EMR Sem Servidor na sua conta. Essas ações podem incorrer em custos para sua Conta da AWS. Ao criar ou editar políticas baseadas em identidade, siga estas diretrizes e recomendações:

- Comece com as políticas AWS gerenciadas e avance para as permissões de privilégios mínimos — Para começar a conceder permissões aos seus usuários e cargas de trabalho, use as políticas AWS gerenciadas que concedem permissões para muitos casos de uso comuns. Eles estão disponíveis no seu Conta da AWS. Recomendamos que você reduza ainda mais as permissões definindo políticas gerenciadas pelo AWS cliente que sejam específicas para seus casos de uso. Para saber mais, consulte [Políticas gerenciadas pela AWS](#) ou [Políticas gerenciadas pela AWS para funções de trabalho](#) no Guia do usuário do IAM.
- Aplique permissões de privilégio mínimo: ao definir permissões com as políticas do IAM, conceda apenas as permissões necessárias para executar uma tarefa. Você faz isso definindo as ações que podem ser executadas em recursos específicos sob condições específicas, também conhecidas como permissões de privilégio mínimo. Para saber mais sobre como usar o IAM para aplicar permissões, consulte [Políticas e permissões no IAM](#) no Guia do usuário do IAM.
- Use condições nas políticas do IAM para restringir ainda mais o acesso: é possível adicionar uma condição às políticas para limitar o acesso a ações e recursos. Por exemplo, é possível escrever uma condição de política para especificar que todas as solicitações devem ser enviadas usando SSL. Você também pode usar condições para conceder acesso às ações de serviço se elas forem usadas por meio de uma ação específica AWS service (Serviço da AWS), como CloudFormation. Para saber mais, consulte [Elementos da política JSON do IAM: condição](#) no Guia do usuário do IAM.
- Use o IAM Access Analyzer para validar suas políticas do IAM a fim de garantir permissões seguras e funcionais: o IAM Access Analyzer valida as políticas novas e existentes para que elas sigam a linguagem de política do IAM (JSON) e as práticas recomendadas do IAM. O IAM Access Analyzer oferece mais de cem verificações de política e recomendações práticas para ajudar

a criar políticas seguras e funcionais. Para saber mais, consulte [Validação de políticas do IAM Access Analyzer](#) no Guia do Usuário do IAM.

- Exigir autenticação multifator (MFA) — Se você tiver um cenário que exija usuários do IAM ou um usuário root, ative Conta da AWS a MFA para obter segurança adicional. Para exigir MFA quando as operações de API forem chamadas, adicione condições de MFA às suas políticas. Para saber mais, consulte [Configuração de acesso à API protegido por MFA](#) no Guia do Usuário do IAM.

Para saber mais sobre as práticas recomendadas do IAM, consulte [Práticas recomendadas de segurança no IAM](#) no Guia do usuário do IAM.

Permitir que os usuários acessem suas próprias permissões

Este exemplo mostra como criar uma política que permita que os usuários do IAM visualizem as políticas gerenciadas e em linha anexadas a sua identidade de usuário. Essa política inclui permissões para concluir essa ação no console ou programaticamente usando a API AWS CLI ou AWS .

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "ViewOwnUserInfo",
      "Effect": "Allow",
      "Action": [
        "iam:GetUserPolicy",
        "iam:ListGroupsWithUser",
        "iam:ListAttachedUserPolicies",
        "iam:ListUserPolicies",
        "iam:GetUser"
      ],
      "Resource": ["arn:aws:iam::*:user/${aws:username}"]
    },
    {
      "Sid": "NavigateInConsole",
      "Effect": "Allow",
      "Action": [
        "iam:GetGroupPolicy",
        "iam:GetPolicyVersion",
        "iam:GetPolicy",
        "iam:ListAttachedGroupPolicies",
        "iam:ListGroupPolicies",

```

```

        "iam:ListPolicyVersions",
        "iam:ListPolicies",
        "iam:ListUsers"
    ],
    "Resource": "*"
}
]
}

```

Atualizações sem servidor do Amazon EMR para políticas gerenciadas AWS

Acesse detalhes sobre as atualizações das políticas AWS gerenciadas do Amazon EMR Serverless desde que esse serviço começou a rastrear essas alterações. Para receber alertas automáticos sobre mudanças nesta página, assine o feed RSS na página [Histórico de documentos](#) do Amazon EMR Sem Servidor.

Alteração	Descrição	Data
Amazon EMRServerless ServiceRolePolicy — Atualização de uma política existente	O Amazon EMR Serverless adicionou o novo e SidCloudWatchPolicyStatement à EC2PolicyStatement política da Amazon. EMRServerless ServiceRolePolicy	25 de janeiro de 2024
Amazon EMRServerless ServiceRolePolicy — Atualização de uma política existente	O Amazon EMR Sem Servidor adicionou novas permissões para permitir que o Amazon EMR Sem Servidor publique métricas de conta agregadas para uso de vCPU no namespace "AWS/Usage" .	20 de abril de 2023

Alteração	Descrição	Data
O Amazon EMR Sem Servidor passou a monitorar alterações	O Amazon EMR Serverless começou a monitorar as alterações em suas políticas gerenciadas. AWS	20 de abril de 2023

Solução de problemas de identidade e acesso da Amazon EMR Sem Servidor

Use as informações a seguir para ajudar a diagnosticar e corrigir problemas comuns que você pode encontrar ao trabalhar com o Amazon EMR Sem Servidor e o IAM.

Tópicos

- [Não tenho autorização para executar uma ação no Amazon EMR Sem Servidor](#)
- [Não estou autorizado a realizar iam: PassRole](#)
- [Quero permitir que pessoas fora da minha AWS conta acessem meus recursos sem servidor do Amazon EMR](#)
- [Não consigo abrir o servidor de UI/Spark histórico ao vivo do EMR Studio para depurar meu trabalho ou ocorre um erro de API quando tento obter registros usando get-dashboard-for-job-run](#)

Não tenho autorização para executar uma ação no Amazon EMR Sem Servidor

Se isso Console de gerenciamento da AWS indicar que você não está autorizado a realizar uma ação, entre em contato com o administrador para obter ajuda. O administrador é a pessoa que forneceu o seu nome de usuário e senha.

O exemplo de erro a seguir ocorre quando o usuário mateojackson tenta usar o console para acessar detalhes sobre um recurso fictício do *my-example-widget*, mas não tem as permissões fictícias do `emr-serverless:GetWidget`.

```
User: arn:aws:iam::123456789012:user/mateojackson is not authorized to perform: emr-serverless:GetWidget on resource: my-example-widget
```

Neste caso, Mateo pede ao administrador para atualizar suas políticas para permitir a ele o acesso ao recurso *my-example-widget* usando a ação `emr-serverless:GetWidget`.

Não estou autorizado a realizar iam: PassRole

Caso receba uma mensagem de erro informando que você não tem autorização para executar a ação `iam:PassRole`, as políticas deverão ser atualizadas para permitir a transmissão de um perfil ao Amazon EMR Sem Servidor.

Alguns Serviços da AWS permitem que você passe uma função existente para esse serviço em vez de criar uma nova função de serviço ou uma função vinculada ao serviço. Para fazer isso, é preciso ter permissões para passar o perfil para o serviço.

O erro exemplificado a seguir ocorre quando uma usuária do IAM chamada `marymajor` tenta usar o console para executar uma ação no Amazon EMR Sem Servidor. No entanto, a ação exige que o serviço tenha permissões concedidas por um perfil de serviço. Mary não tem permissões para passar o perfil para o serviço.

```
User: arn:aws:iam::123456789012:user/marymajor is not authorized to perform:
iam:PassRole
```

Nesse caso, as políticas de Mary devem ser atualizadas para permitir que ela realize a ação `iam:PassRole`.

Se precisar de ajuda, entre em contato com seu AWS administrador. Seu administrador é a pessoa que forneceu suas credenciais de login.

Quero permitir que pessoas fora da minha AWS conta acessem meus recursos sem servidor do Amazon EMR

É possível criar um perfil que os usuários de outras contas ou pessoas fora da organização podem usar para acessar seus recursos. É possível especificar quem é confiável para assumir o perfil. Para serviços que oferecem suporte a políticas baseadas em recursos ou listas de controle de acesso (ACLs), você pode usar essas políticas para conceder às pessoas acesso aos seus recursos.

Para saber mais, consulte:

- Para saber se o Amazon EMR Sem Servidor é compatível com esses recursos, consulte [Identity and Access Management \(IAM\) no Amazon EMR Sem Servidor](#).
- Para saber como fornecer acesso aos seus recursos em todos os Contas da AWS que você possui, consulte [Como fornecer acesso a um usuário do IAM em outro Conta da AWS que você possui](#) no Guia do usuário do IAM.

- Para saber como fornecer acesso aos seus recursos a terceiros Contas da AWS, consulte [Como fornecer acesso Contas da AWS a terceiros](#) no Guia do usuário do IAM.
- Para saber como conceder acesso por meio da federação de identidades, consulte [Conceder acesso a usuários autenticados externamente \(federação de identidades\)](#) no Guia do usuário do IAM.
- Para saber a diferença entre perfis e políticas baseadas em recurso para acesso entre contas, consulte [Acesso a recursos entre contas no IAM](#) no Guia do usuário do IAM.

Não consigo abrir o servidor de UI/Spark histórico ao vivo do EMR Studio para depurar meu trabalho ou ocorre um erro de API quando tento obter registros usando **get-dashboard-for-job-run**

Se você usa o armazenamento gerenciado do EMR Sem Servidor para registro em log, e sua aplicação do EMR Sem Servidor está em uma sub-rede privada com endpoints de VPC para o Amazon S3 e você anexa uma política de endpoint para controlar o acesso, adicione as permissões mencionadas em [Registro em log para o EMR Sem Servidor com armazenamento gerenciado](#) na sua política de VPC ao endpoint de gateway S3 para o EMR Sem Servidor armazenar e veicular logs da aplicação.

Propagação de identidades confiáveis

Com as versões 7.8.0 e superiores do Amazon EMR, você pode propagar identidades de usuários do AWS IAM Identity Center para cargas de trabalho interativas com o EMR Serverless por meio do Apache Livy Endpoint. As workloads interativas do Apache Livy propagarão ainda mais a identidade fornecida para serviços downstream, como Amazon S3, Lake Formation e Amazon Redshift, permitindo o acesso seguro aos dados por meio da identidade do usuário nesses serviços. As seções a seguir fornecem uma visão geral conceitual, os pré-requisitos e as etapas necessárias para iniciar e propagar a identidade para workloads interativas com o EMR Sem Servidor por meio do Endpoint Apache Livy.

Visão geral do

O [IAM Identity Center](#) é a abordagem recomendada para autenticação e autorização da força de trabalho em AWS organizações de qualquer tamanho e tipo. Com o Identity Center, crie e gerencie identidades de usuários ou conecte sua fonte de identidade existente, incluindo Microsoft Active

Directory, Okta, Ping Identity JumpCloud, Google Workspace e Microsoft Entra ID (antigo Azure AD).
AWS

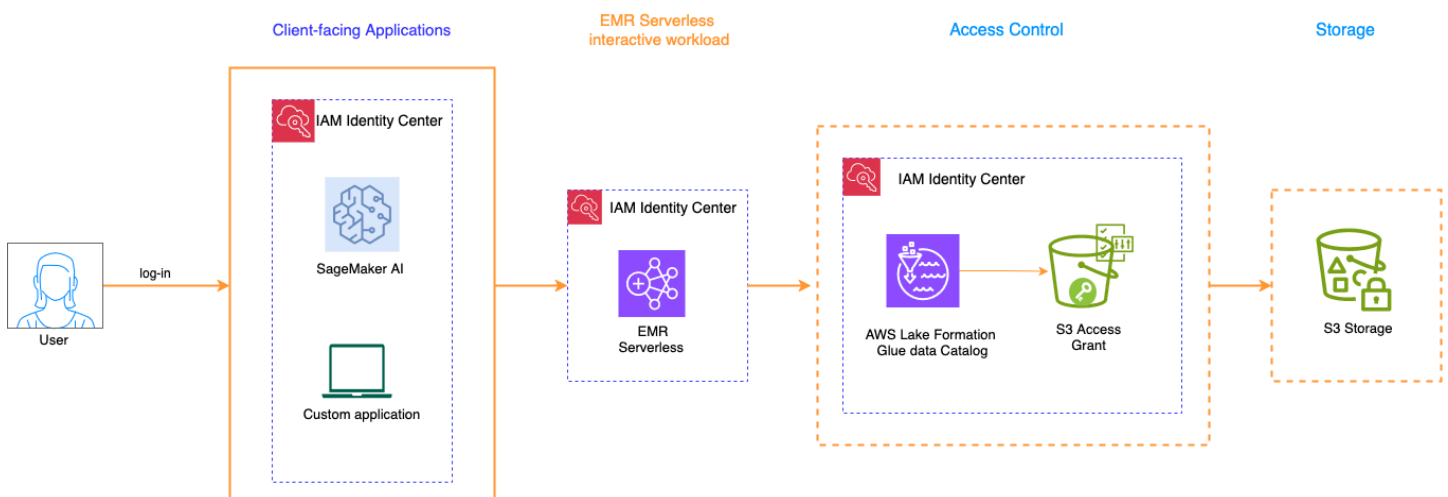
A [propagação confiável de identidade](#) é um recurso AWS do IAM Identity Center que os administradores de AWS serviços conectados podem usar para conceder e auditar o acesso aos dados do serviço. O acesso a esses dados é baseado em atributos do usuário, como associações de grupo. Configurar a propagação de identidade confiável requer colaboração entre os administradores dos AWS serviços conectados e os administradores do IAM Identity Center. Para obter mais informações, consulte [Pré-requisitos e considerações](#), no Guia do usuário do Centro de Identidade do IAM.

Recursos e benefícios

A integração do Endpoint Apache Livy do EMR Sem Servidor com a [Propagação de identidade confiável](#) do Centro de Identidade do IAM oferece os seguintes benefícios:

- A capacidade de impor a autorização em nível de tabela com identidades do Identity Center nas tabelas do catálogo de dados AWS Glue gerenciado pela Lake AWS Formation.
- A capacidade de aplicar a autorização com identidades do Centro de Identidade em clusters do Amazon Redshift.
- Permite o rastreamento ponta a ponta das ações do usuário para fins de auditoria.
- A capacidade de aplicar a autorização no nível de prefixo do Amazon S3 com as identidades do Centro de Identidade em prefixos do S3 gerenciados pelo S3 Access Grants.

Como funciona



Exemplos de casos de uso

Preparação de dados e engenharia de atributos

Cientistas de dados de várias equipes de pesquisa colaboram em projetos complexos usando uma plataforma de dados unificada. Eles se conectam à SageMaker IA usando suas credenciais corporativas, obtendo acesso imediato a um vasto lago de dados compartilhado que abrange várias AWS contas. À medida que iniciam a engenharia de atributos para novos modelos de machine learning, as sessões do Spark lançadas por meio do EMR Sem Servidor reforçam as políticas de segurança em nível de coluna e linha do Lake Formation com base em suas identidades propagadas. Os cientistas podem preparar dados de forma eficiente e criar recursos usando ferramentas familiares, enquanto as equipes de conformidade têm a garantia de que cada interação de dados será automaticamente rastreada e auditada. Esse ambiente seguro e colaborativo acelera os pipelines de pesquisa enquanto mantém os rígidos padrões de proteção de dados exigidos em setores regulamentados.

Introdução à propagação de identidade confiável

[Esta seção ajuda você a configurar o aplicativo EMR-Serverless com o Apache Livy Endpoint para integrar-se ao AWS IAM Identity Center e permitir a propagação confiável de identidades.](#)

Pré-requisitos

- Uma instância do Identity Center na AWS região em que você deseja criar uma propagação de identidade confiável habilitada para EMR Serverless Apache Livy Endpoint. Uma instância do Identity Center só pode existir em uma única região para uma AWS conta. Consulte [Habilitar o IAM Identity Center](#) e [provisione os usuários e grupos da sua fonte de identidades no IAM Identity Center](#).
- Habilite a propagação de identidade confiável para serviços downstream, como Lake Formation, Concessão de acesso do S3 ou cluster do Amazon Redshift, com os quais a carga de trabalho interativa interage para acessar dados.

Permissões para criar o aplicativo EMR Serverless habilitado para propagação de identidade confiável

Além das [permissões básicas necessárias para acessar o EMR Sem Servidor](#), você deve configurar permissões adicionais para sua identidade ou perfil do IAM usado para criar uma aplicação do EMR Sem Servidor habilitada para Propagação de identidade confiável. Para propagação de identidade

confiável, o EMR creates/bootstraps Serverless é um aplicativo de centro de identidade gerenciado por um único serviço em sua conta, que o serviço utiliza para validação de identidade e propagação de identidade para downstream.


```
"sso:DescribeInstance",  
"sso:CreateApplication",  
"sso>DeleteApplication",  
"sso:PutApplicationAuthenticationMethod",  
"sso:PutApplicationAssignmentConfiguration",  
"sso:PutApplicationGrant",  
"sso:PutApplicationAccessScope"
```

- `sso:DescribeInstance`— Concede permissão para descrever e validar o IAM Identity Center InstanceARN que você especifica no parâmetro. `identity-center-configuration`
- `sso:CreateApplication`— Concede permissão para criar um aplicativo IAM Identity Center gerenciado pelo EMR Serverless, que é usado para ações. `trusted-identity-propatgion`
- `sso>DeleteApplication`: concede permissão para limpar uma aplicação do Centro de Identidade do IAM gerenciada pelo EMR Sem Servidor
- `sso:PutApplicationAuthenticationMethod`: concede permissão para colocar o método de autenticação na aplicação do Centro de Identidade do IAM gerenciada pelo EMR Sem Servidor, que permite que a entidade de serviço `emr-serverless` interaja com a aplicação do Centro de Identidade do IAM.
- `sso:PutApplicationAssignmentConfiguration`— Concede permissão para definir a configuração `ser-assignment-not-required` “U” no aplicativo IAM Identity Center.
- `sso:PutApplicationGrant`: Concede permissão para aplicar concessões de troca de token, `introspectToken`, `refreshToken` e `revokeToken` em uma aplicação do Centro de Identidade do IAM.
- `sso:PutApplicationAccessScope`: concede permissão para aplicar o escopo downstream habilitado para propagação de identidade confiável à aplicação do Centro de Identidade do IAM. Aplicamos os escopos `redshift:connect`, `lakeformation:query` e `s3:read_write` para habilitar esses serviços. `trusted-identity-propagation`

Criar uma aplicação do EMR Sem Servidor habilitada para propagação de identidades confiáveis

Você deve especificar o campo `-identity-center-configuration` com `identityCenterInstanceArn` para permitir a propagação de identidades confiáveis na aplicação.

Use o comando de exemplo a seguir para criar uma aplicação do EMR Sem Servidor com a propagação de identidade confiável habilitada.

 Note

Você também deve especificar `--interactive-configuration '{"livyEndpointEnabled":true}'` porque a propagação de identidades confiáveis está habilitada apenas para o Endpoint Apache Livy.

```
aws emr-serverless create-application \  
  --release-label emr-7.8.0 \  
  --type "SPARK" \  
  --identity-center-configuration '{"identityCenterInstanceArn" :  
"arn:aws:sso:::instance/ssoins-123456789"}' \  
  --interactive-configuration '{"livyEndpointEnabled":true}'
```

- `identity-center-configuration`: (opcional) Habilita a propagação de identidade confiável do Centro de Identidade, se especificado.
- `identityCenterInstanceArn`: (obrigatório) o ARN da instância do Centro de Identidade.

Caso você não tenha as permissões necessárias do Centro de Identidade (mencionadas anteriormente), primeiro crie a aplicação do EMR Sem Servidor sem propagação de identidades confiáveis (por exemplo, não especifique o parâmetro `--identity-center-configuration`) e, posteriormente, solicite ao administrador do Centro de Identidade que habilite a propagação de identidades confiáveis invocando a API `update-application`. Veja o exemplo abaixo:

```
aws emr-serverless update-application \  
  --application-id applicationId \  
  --identity-center-configuration '{"identityCenterInstanceArn" :  
"arn:aws:sso:::instance/ssoins-123456789"}'
```

O EMR Sem Servidor cria uma aplicação do Centro de Identidade gerenciada por serviço na sua conta, que o serviço utiliza para validações de identidade e propagação de identidade para serviços downstream. O aplicativo Identity Center gerenciado criado pelo EMR Serverless é compartilhado entre todos os aplicativos `trusted-identity-propagation` EMR Serverless habilitados em sua conta.

Note

Não modifique manualmente as configurações na aplicação do Centro de Identidade gerenciada. Qualquer alteração pode afetar todos os aplicativos EMR Serverless trusted-identity-propagation habilitados em sua conta.

Permissões do perfil de execução de trabalho para propagar identidade

Como o EMR-Serverless utiliza job-execution-role credenciais aprimoradas de identidade para propagar a identidade para AWS serviços posteriores, a política de confiança da Job Execution Role deve ter `sts:SetContext` permissão adicional para aprimorar a credencial da função de execução de tarefas com identidade para permitir serviços posteriores, como S3 access-grant trusted-identity-propagation, Lake Formation ou Amazon Redshift. Para saber mais sobre como criar um perfil, consulte [Criar um perfil de runtime de trabalho](#).

JSON

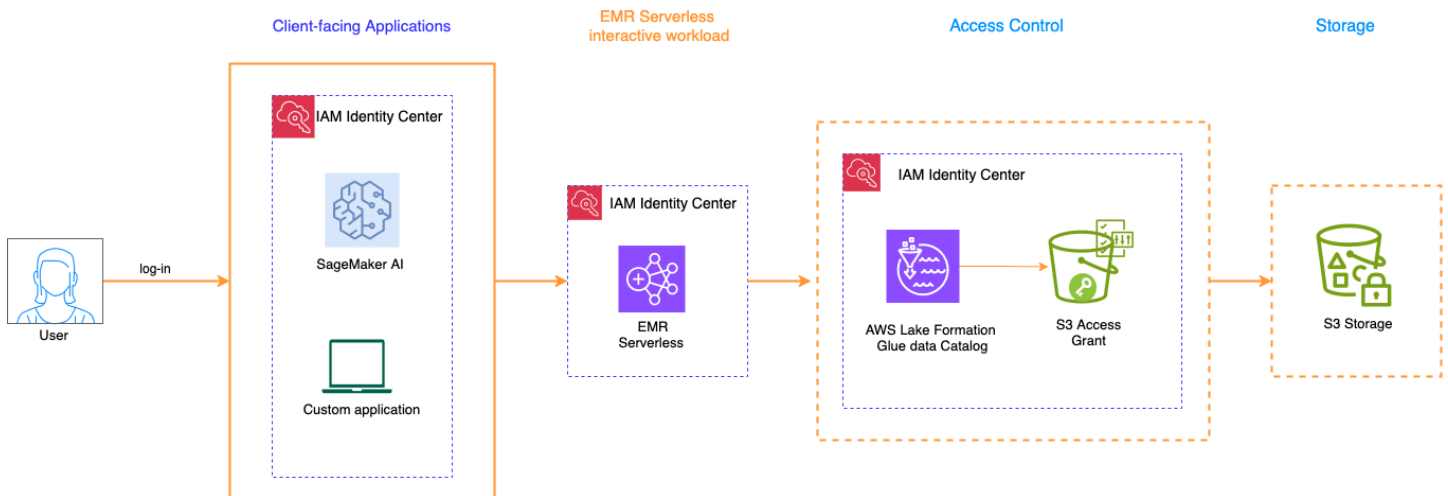
```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "emr-serverless.amazonaws.com"
      },
      "Action": [ "sts:AssumeRole", "sts:SetContext" ]
    }
  ]
}
```

Além disso, JobExecutionRole precisa de permissões para AWS serviços posteriores que o job-run invocaria para buscar dados usando a identidade do usuário. Consulte os links abaixo para configurar a Concessão de acesso do S3, Lake Formation.

- [Usar o Lake Formation com o EMR Sem Servidor](#)
- [Usando as concessões de acesso do Amazon S3 com o EMR Serverless](#)

Propagação de identidade confiável para cargas de trabalho interativas

As etapas para propagar a identidade para cargas de trabalho interativas por meio de um endpoint Apache Livy dependem de seus usuários interagirem com um ambiente de desenvolvimento AWS gerenciado, como seu próprio ambiente de Notebook auto-hospedado como Amazon SageMaker AI aplicativo voltado para o cliente.



AWS ambiente de desenvolvimento gerenciado

O seguinte aplicativo AWS gerenciado voltado para o cliente oferece suporte à propagação confiável de identidade com o endpoint Apache Livy sem servidor EMR:

- [SageMaker Inteligência Artificial da Amazon](#)

Ambiente de notebook auto-hospedado e gerenciado pelo cliente

Para permitir a propagação de identidade confiável para usuários de aplicativos desenvolvidos sob medida, consulte [Acesse os AWS serviços de forma programática usando a propagação de identidade confiável](#) no Blog de Segurança.AWS

Sessões de usuário em segundo plano

As sessões em segundo plano do usuário permitem que os fluxos de análise e aprendizado de máquina de longa duração continuem mesmo depois que o usuário se desconecta da interface do notebook. Esse recurso é implementado por meio da integração do EMR Serverless com o recurso confiável de propagação de identidade do IAM Identity Center. Esta seção explica as opções de configuração e os comportamentos das sessões em segundo plano do usuário.

Note

As sessões em segundo plano do usuário se aplicam às cargas de trabalho do Spark iniciadas por meio de interfaces de notebook, como o Amazon SageMaker Unified Studio. A ativação ou desativação desse recurso afeta somente as novas sessões do Livy; as sessões ativas existentes do Livy não são afetadas.

Configurar sessões de usuário em segundo plano

As sessões de usuário em segundo plano devem ser habilitadas em dois níveis para que funcionem adequadamente:

1. Nível de instância do IAM Identity Center — normalmente configurado pelos administradores do iDC
2. Nível do aplicativo EMR Serverless — configurado pelos administradores do aplicativo EMR Serverless

Habilite sessões em segundo plano do usuário para aplicativos EMR Serverless

Para habilitar sessões em segundo plano do usuário para um aplicativo EMR Serverless, você deve definir o `userBackgroundSessionsEnabled` parâmetro como `true` in `identityCenterConfiguration` ao criar ou atualizar um aplicativo.

Pré-requisitos

- Sua função do IAM usada `create/update` no aplicativo EMR Serverless deve ter a permissão `sso:PutApplicationSessionConfiguration`. Essa permissão permite que o EMR Serverless habilite sessões em segundo plano do usuário no nível do aplicativo iDC gerenciado pelo EMR Serverless.
- Seu aplicativo EMR Serverless deve usar a etiqueta de versão 7.8 ou posterior e ter a Propagação de Identidade Confiável ativada.

Para habilitar sessões em segundo plano do usuário usando o AWS CLI

```
aws emr-serverless create-application \  
  --name "my-analytics-app" \  
  --type "SPARK" \  
  --userBackgroundSessionsEnabled true
```

```
--release-label "emr-7.8.0" \
--identity-center-configuration '{"identityCenterInstanceArn":
"arn:aws:sso:::instance/ssoins-1234567890abcdef", "userBackgroundSessionsEnabled":
true}'
```

Para atualizar um aplicativo existente:

```
aws emr-serverless update-application \
--application-id applicationId \
--identity-center-configuration '{"identityCenterInstanceArn":
"arn:aws:sso:::instance/ssoins-1234567890abcdef", "userBackgroundSessionsEnabled":
true}'
```

Matriz de configuração

A configuração efetiva da sessão em segundo plano do usuário depende da configuração do aplicativo EMR Serverless e das configurações em nível de instância do IAM Identity Center:

Matriz de configuração da sessão em segundo plano do usuário

Centro de identidade do IAM userBackgroundSession ativado	Habilitado para EMR Serverless userBackg roundSessions	Comportamento
Sim	TRUE	Sessões de usuário em segundo plano ativadas
Sim	FALSE	A sessão expira com o logout do usuário
Não	TRUE	O aplicativo creation/update falha com exceção
Não	FALSE	A sessão expira com o logout do usuário

Duração padrão das sessões de usuários em segundo plano

Por padrão, todas as sessões de usuário em segundo plano têm um limite de duração de 7 dias no Centro de Identidade do IAM. Os administradores podem modificar essa duração no console

do Centro de Identidade do IAM. Essa configuração se aplica no nível da instância do Centro de Identidade do IAM, afetando todos os aplicativos do Centro de Identidade do IAM compatíveis nessa instância.

- A duração pode ser definida para qualquer valor, de 15 minutos a 90 dias.
- Essa configuração é definida no console do IAM Identity Center em Configurações → Autenticação → Configurar (seção Trabalhos não interativos)

Note

As sessões do EMR Serverless Livy têm um limite de duração máxima separado de 24 horas. As sessões serão encerradas quando o limite da sessão do Livy ou a duração da sessão em segundo plano do usuário forem atingidos, o que ocorrer primeiro.

Impacto de desativar as sessões de usuário em segundo plano

Quando as sessões em segundo plano do usuário são desativadas no IAM Identity Center:

Sessões existentes do Livy

Continue a ser executado sem interrupção se eles tiverem sido iniciados com as sessões em segundo plano do usuário ativadas. Essas sessões continuarão usando seus tokens de sessão em segundo plano existentes até que terminem naturalmente ou sejam explicitamente interrompidas.

Novas sessões de Livy

Usará o fluxo padrão de propagação de identidade confiável e será encerrado quando o usuário se desconectar ou sua sessão interativa expirar (como ao fechar um notebook Amazon SageMaker Unified Studio JupyterLab).

Alterando a duração das sessões de usuário em segundo plano

Quando a configuração de duração das sessões de usuário em segundo plano é modificada no Centro de Identidade do IAM:

Sessões existentes do Livy

Continue executando com a mesma duração da sessão em segundo plano com a qual eles foram iniciados.

Novas sessões de Livy

Usará a nova duração da sessão para sessões em segundo plano.

Considerações

Condições de encerramento da sessão

Ao usar sessões em segundo plano do usuário, uma sessão do Livy continuará em execução até que uma das seguintes situações ocorra:

- A sessão em segundo plano do usuário expira (com base na configuração do iDC, em até 90 dias)
- A sessão de usuário em segundo plano for revogada manualmente por um administrador.
- A sessão do Livy atinge seu tempo limite de inatividade (padrão: 1 hora após a última instrução executada)
- A sessão Livy atinge sua duração máxima (24 horas)
- O usuário interrompe ou reinicia explicitamente o kernel do notebook

Persistência de dados

Ao usar sessões de usuário em segundo plano:

- Os usuários não podem se reconectar à interface do notebook para visualizar os resultados depois de se desconectarem
- Configure suas instruções do Spark para gravar os resultados no armazenamento persistente (como o Amazon S3) antes que a execução seja concluída

Implicações de custo

- Os trabalhos continuarão sendo executados até a conclusão mesmo depois que os usuários encerrarem a JupyterLab sessão do Amazon SageMaker Unified Studio e incorrerão em cobranças por toda a duração da execução concluída.

- Monitore suas sessões ativas em segundo plano para evitar custos desnecessários de sessões esquecidas ou abandonadas.

Disponibilidade de recursos

As sessões em segundo plano do usuário do EMR Serverless estão disponíveis para:

- Somente motor Spark (o motor Hive não é suportado)
- Somente sessões interativas do Livy (trabalhos em lote e trabalhos de streaming não são suportados)
- Etiquetas de lançamento do EMR Serverless 7.8 e versões posteriores

Considerações sobre a integração sem servidor do Trusted-Identity-Propagation EMR

Considere o seguinte ao usar o IAM Identity Center Trusted-Identity-Propagation com o aplicativo EMR Serverless:

- A propagação de identidade confiável por meio do Centro de Identidade é compatível com o Amazon EMR 7.8.0 e versões superiores, somente com o Apache Spark.
- A Propagação de identidades confiáveis apenas pode ser usada para workloads interativas com o EMR Sem Servidor por meio de um endpoint Apache Livy. Workflows de trabalho interativas por meio do EMR Studio não oferecem suporte ao Propagação de identidades confiáveis
- Os trabalhos em lote e as cargas de trabalho de streaming não oferecem suporte à propagação confiável de identidade
- Controles de acesso refinados usando o AWS Lake Formation que usam o Trusted Identity Propagation estão disponíveis para [cargas de trabalho interativas com o EMR Serverless](#) por meio de um endpoint Apache Livy.
- A propagação de identidade confiável com o Amazon EMR é suportada nas seguintes AWS regiões:
 - af-south-1: África (Cidade do Cabo)
 - ap-east-1: Ásia-Pacífico (Hong Kong)
 - ap-northeast-1 – Ásia-Pacífico (Tóquio)
 - ap-northeast-2: Ásia-Pacífico (Seul)

- Ásia-Pacífico (Osaka): ap-northeast-3
- ap-south-1: Ásia-Pacífico (Mumbai)
- ap-southeast: Ásia-Pacífico (Singapura)
- ap-southeast-2: Ásia-Pacífico (Sydney)
- ap-southeast-3: Ásia-Pacífico (Jacarta)
- ca-central-1: Canadá (Central)
- ca-west-1: Canadá (Calgary)
- eu-central-1 – Europa (Frankfurt)
- eu-north-1: Europa (Estocolmo)
- eu-south-1: Europa (Milão)
- eu-south-2: Europa (Espanha)
- eu-west-1 – Europa (Irlanda)
- eu-west-2: Europa (Londres)
- eu-west-3: Europa (Paris)
- me-central-1: Oriente Médio (Emiratos Árabes Unidos)
- me-south-1: Oriente Médio (Bahrein)
- sa-east-1: América do Sul (São Paulo)
- us-east-1 – Leste dos EUA (Norte da Virgínia)
- us-east-2 – Leste dos EUA (Ohio)
- us-west-1: Oeste dos EUA (Norte da Califórnia)
- us-west-2 – Oeste dos EUA (Oregon)

Usar o Lake Formation com o EMR Sem Servidor

Você pode configurar aplicativos EMR Serverless para usar o Lake Formation com acesso total à tabela ou controle de acesso refinado. Para obter detalhes sobre os recursos compatíveis em cada modo de acesso, consulte a tabela a seguir.

Disponibilidade de recursos

Recurso	Disponível a partir de
Operações de leitura (SELECT, DESCRIBE) para tabelas Hive e Iceberg	EMR 7.2+
Visualizações em vários dialetos	EMR 7.6+
Operações de leitura (SELECT, DESCRIBE) para tabelas Delta Lake e Hudi	EMR 7.6+
Acesso completo à mesa para Hive, Iceberg	EMR 7.9+
Acesso completo à mesa para Delta Lake	EMR 7.11+
Operações de gravação (DDL, DML) para tabelas Hive, Iceberg e Delta Lake	EMR 7.12+
Acesso total à mesa para Hudi	EMR 7.12+

Acesso completo à tabela do Lake Formation para o EMR Serverless

Com as versões 7.8.0 e superiores do Amazon EMR, você pode aproveitar o Lake AWS Formation com o Glue Data Catalog, onde a função de tempo de execução do trabalho tem permissões completas de tabela sem as limitações do controle de acesso refinado. Esse recurso permite que você leia e grave em tabelas protegidas pelo Lake Formation a partir de trabalhos Spark em lote e interativos do EMR Sem Servidor. Consulte as seções a seguir para saber mais sobre o Lake Formation e como usá-lo com o EMR Sem Servidor.

Como usar o Lake Formation com acesso total à tabela

Você pode acessar as tabelas do catálogo do Glue Data protegidas pelo AWS Lake Formation a partir de trabalhos do EMR Serverless Spark ou sessões interativas em que a função de tempo de execução do trabalho tem acesso total à tabela. Você não precisa habilitar o AWS Lake Formation no aplicativo EMR Serverless. Quando um trabalho do Spark é configurado para Full Table Access (FTA), as credenciais do AWS Lake Formation são usadas para os dados do read/write S3 das

tabelas registradas do AWS Lake Formation, enquanto as credenciais da função de tempo de execução do trabalho são usadas para read/write tabelas não registradas no Lake Formation. AWS

Important

Não ative o AWS Lake Formation para um controle de acesso refinado. Um trabalho não pode executar simultaneamente o Acesso total à tabela (FTA) e o Controle de acesso refinado (FGAC) no mesmo cluster ou aplicação EMR.

Etapa 1: habilitar o acesso total à tabela no Lake Formation

Para usar o modo Full Table Access (FTA), você deve permitir que mecanismos de consulta de terceiros acessem dados sem a validação da tag de sessão do IAM no AWS Lake Formation. Para habilitar, siga as etapas em [Application integration for full table access](#).

Note

Ao acessar tabelas entre contas, o acesso completo à tabela deve ser habilitado nas contas de produtor e consumidor. Da mesma forma, ao acessar tabelas entre regiões, essa configuração deve ser habilitada nas regiões produtora e consumidora.

Etapa 2: configurar permissões do IAM para o perfil de runtime do trabalho

Para acesso de leitura ou gravação aos dados subjacentes, além das permissões do Lake Formation, um perfil de runtime do trabalho precisa da permissão `lakeformation:GetDataAccess` do IAM. Com essa permissão, o Lake Formation concede a solicitação de credenciais temporárias para acessar os dados.

Confira a seguir um exemplo de política de como fornecer permissões do IAM para acesso a um script no S3, upload de logs no Amazon S3, permissões da API do AWS Glue e permissão para acessar o Lake Formation.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
```

```
{
  "Sid": "ScriptAccess",
  "Effect": "Allow",
  "Action": [
    "s3:GetObject"
  ],
  "Resource": [
    "arn:aws:s3:::*.amzn-s3-demo-bucket/scripts"
  ]
},
{
  "Sid": "LoggingAccess",
  "Effect": "Allow",
  "Action": [
    "s3:PutObject"
  ],
  "Resource": [
    "arn:aws:s3:::amzn-s3-demo-bucket/logs/*"
  ]
},
{
  "Sid": "GlueCatalogAccess",
  "Effect": "Allow",
  "Action": [
    "glue:Get*",
    "glue:Create*",
    "glue:Update*"
  ],
  "Resource": [
    "*"
  ]
},
{
  "Sid": "LakeFormationAccess",
  "Effect": "Allow",
  "Action": [
    "lakeformation:GetDataAccess"
  ],
  "Resource": [
    "*"
  ]
}
]
```

```
}
```

Etapa 2.1: configurar permissões do Lake Formation

- Os trabalhos Spark que leem dados do S3 necessitam de permissão SELECT do Lake Formation.
- O Spark faz com que write/delete os dados no S3 exijam a permissão Lake Formation ALL (SUPER).
- Os trabalhos do Spark que interagem com o Catálogo de dados do Glue exigem as permissões DESCRIBE, ALTER, DROP, conforme apropriado.

Para obter mais informações, consulte [Conceder permissões em recursos do Catálogo de dados](#).

Etapa 3: inicializar uma sessão do Spark para acesso total à tabela usando o Lake Formation

Pré-requisitos

AWS O Glue Data Catalog deve ser configurado como um metastore para acessar as tabelas do Lake Formation.

Defina as seguintes configurações para configurar o catálogo do Glue como um metastore:

```
--conf spark.sql.catalogImplementation=hive
--conf
  spark.hive.metastore.client.factory.class=com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalog
```

Para obter mais informações sobre como habilitar o Catálogo de dados para o EMR Sem Servidor, consulte [Configuração do metastore para o EMR Sem Servidor](#).

Para acessar tabelas registradas no AWS Lake Formation, as seguintes configurações precisam ser definidas durante a inicialização do Spark para configurar o Spark para usar as credenciais do Lake Formation AWS .

Hive

```
--conf
  spark.hadoop.fs.s3.credentialsResolverClass=com.amazonaws.glue.accesscontrol.AWSLakeFormation
--conf spark.hadoop.fs.s3.useDirectoryHeaderAsFolderObject=true
--conf spark.hadoop.fs.s3.folderObject.autoAction.disabled=true
--conf spark.sql.catalog.skipLocationValidationOnCreateTable.enabled=true
--conf spark.sql.catalog.createDirectoryAfterTable.enabled=true
```

```
--conf spark.sql.catalog.dropDirectoryBeforeTable.enabled=true
```

Iceberg

```
--conf spark.sql.catalog.spark_catalog=org.apache.iceberg.spark.SparkSessionCatalog
--conf spark.sql.catalog.spark_catalog.warehouse=S3_DATA_LOCATION
--conf spark.sql.catalog.spark_catalog.client.region=REGION
--conf spark.sql.catalog.spark_catalog.type=glue
--conf spark.sql.catalog.spark_catalog.glue.account-id=ACCOUNT_ID
--conf spark.sql.catalog.spark_catalog.glue.lakeformation-enabled=true
--conf spark.sql.catalog.dropDirectoryBeforeTable.enabled=true
```

Delta Lake

```
--conf
  spark.hadoop.fs.s3.credentialsResolverClass=com.amazonaws.glue.accesscontrol.AWSLakeFormation
--conf spark.hadoop.fs.s3.useDirectoryHeaderAsFolderObject=true
--conf spark.hadoop.fs.s3.folderObject.autoAction.disabled=true
--conf spark.sql.catalog.skipLocationValidationOnCreateTable.enabled=true
--conf spark.sql.catalog.createDirectoryAfterTable.enabled=true
--conf spark.sql.catalog.dropDirectoryBeforeTable.enabled=true
```

Hudi

```
--conf
  spark.hadoop.fs.s3.credentialsResolverClass=com.amazonaws.glue.accesscontrol.AWSLakeFormation
--conf spark.hadoop.fs.s3.useDirectoryHeaderAsFolderObject=true
--conf spark.hadoop.fs.s3.folderObject.autoAction.disabled=true
--conf spark.sql.catalog.skipLocationValidationOnCreateTable.enabled=true
--conf spark.sql.catalog.createDirectoryAfterTable.enabled=true
--conf spark.sql.catalog.dropDirectoryBeforeTable.enabled=true
--conf spark.jars=/usr/lib/hudi/hudi-spark-bundle.jar
--conf spark.sql.extensions=org.apache.spark.sql.hudi.HoodieSparkSessionExtension
--conf
  spark.sql.catalog.spark_catalog=org.apache.spark.sql.hudi.catalog.HoodieCatalog
--conf spark.serializer=org.apache.spark.serializer.KryoSerializer
```

- `spark.hadoop.fs.s3.credentialsResolverClass=com.amazonaws.glue.accesscontrol.AWSLakeFormation`
Configure o EMR Filesystem (EMRFS) ou o EMR S3A para usar as credenciais do Lake Formation S3 para tabelas registradas do Lake AWS Formation. Se a tabela não estiver registrada, use as credenciais do perfil de runtime do trabalho.

- `spark.hadoop.fs.s3.useDirectoryHeaderAsFolderObject=true` e `spark.hadoop.fs.s3.folderObject.autoAction.disabled=true`: configure o EMRFS para usar o cabeçalho do tipo de conteúdo `application/x-directory` em vez do sufixo `$folder$` ao criar pastas do S3. Isso é necessário ao ler tabelas do Lake Formation, pois as credenciais do Lake Formation não permitem a leitura de pastas de tabelas com o sufixo `$folder$`.
- `spark.sql.catalog.skipLocationValidationOnCreateTable.enabled=true`: configure o Spark para ignorar a validação da ausência do local da tabela antes da criação. Isso é necessário para tabelas registradas do Lake Formation, pois as credenciais do Lake Formation para verificar o local vazio só estão disponíveis após a criação da tabela do Catálogo de dados do Glue. Sem essa configuração, as credenciais do perfil de runtime do trabalho validarão o local da tabela vazia.
- `spark.sql.catalog.createDirectoryAfterTable.enabled=true`: configure o Spark para criar a pasta do Amazon S3 após a criação da tabela no metastore do Hive. Isso é necessário para tabelas registradas no Lake Formation, pois as credenciais do Lake Formation para criar a pasta do S3 só ficam disponíveis após a criação da tabela no Catálogo de dados do Glue.
- `spark.sql.catalog.dropDirectoryBeforeTable.enabled=true`: configure o Spark para encerrar a pasta do S3 antes da exclusão da tabela no metastore do Hive. Isso é necessário para tabelas registradas do Lake Formation, pois as credenciais do Lake Formation para encerrar a pasta do S3 não estão disponíveis após a exclusão da tabela do Catálogo de dados do Glue.
- `spark.sql.catalog.<catalog>.glue.lakeformation-enabled=true`: Configure o catálogo do Iceberg para usar as credenciais do AWS Lake Formation S3 para tabelas registradas do Lake Formation. Se a tabela não estiver registrada, use as credenciais de ambiente padrão.

Configurar o modo de acesso total à tabela no SageMaker Unified Studio

Para acessar as tabelas registradas do Lake Formation a partir de sessões interativas do Spark em JupyterLab notebooks, use o modo de permissão de compatibilidade. Use o comando mágico `%configure` para definir sua configuração do Spark. Escolha a configuração com base no seu tipo de tabela:

For Hive tables

```
%%configure -f
{
  "conf": {
    "spark.hadoop.fs.s3.credentialsResolverClass":
    "com.amazonaws.glue.accesscontrol.AWSLakeFormationCredentialResolver",
```

```

    "spark.hadoop.fs.s3.useDirectoryHeaderAsFolderObject": true,
    "spark.hadoop.fs.s3.folderObject.autoAction.disabled": true,
    "spark.sql.catalog.skipLocationValidationOnCreateTable.enabled": true,
    "spark.sql.catalog.createDirectoryAfterTable.enabled": true,
    "spark.sql.catalog.dropDirectoryBeforeTable.enabled": true
  }
}

```

For Iceberg tables

```

%%configure -f
{
  "conf": {
    "spark.sql.catalog.spark_catalog":
"org.apache.iceberg.spark.SparkSessionCatalog",
    "spark.sql.catalog.spark_catalog.warehouse": "S3_DATA_LOCATION",
    "spark.sql.catalog.spark_catalog.client.region": "REGION",
    "spark.sql.catalog.spark_catalog.type": "glue",
    "spark.sql.catalog.spark_catalog.glue.account-id": "ACCOUNT_ID",
    "spark.sql.catalog.spark_catalog.glue.lakeformation-enabled": "true",
    "spark.sql.catalog.dropDirectoryBeforeTable.enabled": "true",
  }
}

```

For Delta Lake tables

```

%%configure -f
{
  "conf": {
    "spark.hadoop.fs.s3.credentialsResolverClass":
"com.amazonaws.glue.accesscontrol.AWSLakeFormationCredentialResolver",
    "spark.hadoop.fs.s3.useDirectoryHeaderAsFolderObject": true,
    "spark.hadoop.fs.s3.folderObject.autoAction.disabled": true,
    "spark.sql.catalog.skipLocationValidationOnCreateTable.enabled": true,
    "spark.sql.catalog.createDirectoryAfterTable.enabled": true,
    "spark.sql.catalog.dropDirectoryBeforeTable.enabled": true
  }
}

```

For Hudi tables

```

%%configure -f
{

```

```
"conf": {
  "spark.hadoop.fs.s3.credentialsResolverClass":
"com.amazonaws.glue.accesscontrol.AWSLakeFormationCredentialResolver",
  "spark.hadoop.fs.s3.useDirectoryHeaderAsFolderObject": true,
  "spark.hadoop.fs.s3.folderObject.autoAction.disabled": true,
  "spark.sql.catalog.skipLocationValidationOnCreateTable.enabled": true,
  "spark.sql.catalog.createDirectoryAfterTable.enabled": true,
  "spark.sql.catalog.dropDirectoryBeforeTable.enabled": true,
  "spark.jars": "/usr/lib/hudi/hudi-spark-bundle.jar",
  "spark.sql.extensions":
"org.apache.spark.sql.hudi.HoodieSparkSessionExtension",
  "spark.sql.catalog.spark_catalog":
"org.apache.spark.sql.hudi.catalog.HoodieCatalog",
  "spark.serializer": "org.apache.spark.serializer.KryoSerializer"
}
```

Substitua os espaços reservados:

- S3_DATA_LOCATION: caminho do bucket do S3
- REGION: AWS região (por exemplo, us-east-1)
- ACCOUNT_ID: ID AWS da sua conta

Note

Você precisa definir essas configurações antes de executar qualquer operação do Spark em seu caderno.

Operações com Suporte

Essas operações usarão as credenciais do AWS Lake Formation para acessar os dados da tabela.

- CRIAR TABELA
- ALTER TABLE
- INSERT INTO
- INSERT OVERWRITE
- UPDATE

- MERGE INTO
- DELETE FROM
- ANALYZE TABLE
- REPAIR TABLE
- DESCARTAR TABELA
- Consultas de fontes de dados do Spark
- Gravações na fonte de dados do Spark

Note

As operações não listadas acima continuarão usando as permissões do IAM para acessar os dados da tabela.

Considerações

- Se uma tabela do Hive for criada usando um trabalho que não tem acesso total à tabela habilitado e nenhum registro for inserido, as leituras ou gravações subsequentes de um trabalho com acesso total à tabela falharão. Isso ocorre porque o EMR Spark sem acesso total à tabela adiciona o sufixo `$folder$` ao nome da pasta da tabela. Para resolver isso, você tem as seguintes opções:
 - Insira pelo menos uma linha na tabela diretamente de um trabalho no qual o FTA não esteja habilitado.
 - Configure o trabalho no qual o FTA não está habilitado para não usar o sufixo `$folder$` no nome da pasta no S3. É possível fazer isso definindo a configuração `spark.hadoop.fs.s3.useDirectoryHeaderAsFolderObject=true` do Spark.
 - Crie uma pasta S3 no local da tabela `s3://path/to/table/table_name` usando o console do AWS S3 ou a CLI do AWS S3.
- O acesso total à tabela é compatível com o sistema de arquivos do EMR (EMRFS) a partir da versão 7.8.0 do Amazon EMR e com o sistema de arquivos S3A a partir da versão 7.10.0 do Amazon EMR.
- O Full Table Access é compatível com as tabelas Hive, Iceberg, Delta e Hudi.
- Considerações sobre o Hudi FTA Write Support:
 - As gravações do Hudi FTA exigem o uso `HoodieCredentialedHadoopStorage` para venda de credenciais durante a execução do trabalho.

Defina a seguinte configuração ao executar trabalhos do Hudi:

```
hoodie.storage.class=org.apache.spark.sql.hudi.storage.HoodieCredentialedHadoop
```

- O suporte de gravação do Full Table Access (FTA) para Hudi está disponível a partir da versão 7.12 do Amazon EMR.
- Atualmente, o suporte de gravação do Hudi FTA funciona apenas com as configurações padrão do Hudi. As configurações Hudi personalizadas ou não padrão podem não ser totalmente compatíveis e podem resultar em um comportamento inesperado.
- O agrupamento para tabelas Hudi Merge-On-Read (MOR) não é suportado neste momento no modo de gravação FTA.
- Os trabalhos que mencionam tabelas com regras de controle de acesso de alta granularidade (FGAC) do Lake Formation ou visualizações do catálogo de dados do Glue falharão. Para consultar uma tabela com regras de FGAC ou uma visualização do Catálogo de dados do Glue, você precisará usar o modo FGAC. Você pode ativar o modo FGAC seguindo as etapas descritas na AWS documentação: Usando o [EMR Serverless com o Lake AWS Formation](#) para um controle de acesso refinado.
- O acesso total à tabela não é compatível com Spark Streaming.
- Ao gravar o Spark DataFrame em uma tabela do Lake Formation, somente o modo APPEND é compatível com as tabelas Hive e Iceberg:

```
df.write.mode("append").saveAsTable(table_name)
```
- A criação de tabelas externas requer permissões do IAM.
- Como o Lake Formation armazena temporariamente as credenciais em um trabalho Spark, um trabalho em lote do Spark ou uma sessão interativa em execução no momento pode não refletir as alterações de permissão.
- Você deve usar um perfil definido pelo usuário e não um perfil vinculado ao serviço: [Requisitos do Lake Formation para perfis](#).

Hudi FTA Write Support - Operações suportadas

A tabela a seguir mostra as operações de gravação suportadas para tabelas Hudi Copy-On-Write (COW) e Merge-On-Read (MOR) no modo Full Table Access:

Operações de gravação suportadas pelo Hudi FTA

Tipo de tabela	Operation	Comando SQL Write	Status
VACA	INSERT	INSERT INTO TABLE	Compatível
VACA	INSERT	INSERIR NA TABELA - PARTIÇÃO (estática, dinâmica)	Compatível
VACA	INSERT	INSERT OVERWRITE	Compatível
VACA	INSERT	INSERIR SUBSTITUIÇÃO - PARTIÇÃO (estática, dinâmica)	Compatível
UPDATE	UPDATE	UPDATE TABLE	Compatível
VACA	UPDATE	TABELA DE ATUALIZAÇÃO - Alterar partição	Não suportado
DELETE	DELETE	DELETE FROM TABLE	Compatível
ALTER	ALTER	ALTERAR TABELA - RENOMEAR PARA	Não suportado
VACA	ALTER	ALTERE A TABELA - DEFINA AS	Compatível

Tipo de tabela	Operation	Comando SQL Write	Status
		PROPRIEDADES DA TABELA	
VACA	ALTER	ALTERAR TABELA - DESDEFINIR TBLPROPERTIES	Compatível
VACA	ALTER	ALTERAR TABELA - ALTERAR COLUNA	Compatível
VACA	ALTER	ALTERAR TABELA - ADICIONAR COLUNAS	Compatível
VACA	ALTER	ALTERAR TABELA - ADICIONAR PARTIÇÃO	Compatível
VACA	ALTER	ALTERAR TABELA - SOLTAR PARTIÇÃO	Compatível
VACA	ALTER	ALTERAR TABELA - RECUPERAR PARTIÇÕES	Compatível

Tipo de tabela	Operation	Comando SQL Write	Status
VACA	ALTER	REPARAR PARTIÇÕES DE SINCRONIZAÇÃO DE TABELAS	Compatível
DROP	DROP	DESCARTAR TABELA	Compatível
VACA	DROP	TABELA SUSPENSA - PURGA	Compatível
CREATE	CREATE	CRIAR TABELA - Gerenciado	Compatível
VACA	CREATE	CRIAR TABELA - PARTICIONAR POR	Compatível
VACA	CREATE	CRIAR TABELA SE NÃO EXISTIR	Compatível
VACA	CREATE	CREATE TABLE LIKE	Compatível
VACA	CREATE	CRIAR TABELA COMO SELEÇÃO	Compatível
CREATE	CREATE	CRIAR TABELA com LOCALIZAÇÃO - Tabela externa	Não suportado

Tipo de tabela	Operation	Comando SQL Write	Status
QUADRO DE DADOS (INSERIR)	QUADRO DE DADOS (INSERIR)	saveAsTable.Sobrescrever	Compatível
VACA	QUADRO DE DADOS (INSERIR)	saveAsTable.Anexar	Não suportado
VACA	QUADRO DE DADOS (INSERIR)	saveAsTable.Ignorar	Compatível
VACA	QUADRO DE DADOS (INSERIR)	saveAsTable.ErrorIfExists	Compatível
VACA	QUADRO DE DADOS (INSERIR)	saveAsTable - Tabela externa (Caminho)	Não suportado
VACA	QUADRO DE DADOS (INSERIR)	salvar (caminho) - DF v1	Não suportado
MAIS	INSERT	INSERT INTO TABLE	Compatível
MAIS	INSERT	INSERIR NA TABELA - PARTIÇÃO (estática, dinâmica)	Compatível
MAIS	INSERT	INSERT OVERWRITE	Compatível

Tipo de tabela	Operation	Comando SQL Write	Status
MAIS	INSERT	INSERIR SUBSTITUIÇÃO - PARTIÇÃO (estática, dinâmica)	Compatível
UPDATE	UPDATE	UPDATE TABLE	Compatível
MAIS	UPDATE	TABELA DE ATUALIZAÇÃO - Alterar partição	Não suportado
DELETE	DELETE	DELETE FROM TABLE	Compatível
ALTER	ALTER	ALTERAR TABELA - RENOMEAR PARA	Não suportado
MAIS	ALTER	ALTERE A TABELA - DEFINA AS PROPRIEDA DES DA TABELA	Compatível
MAIS	ALTER	ALTERAR TABELA - DESDEFINIR TBLPROPER TIES	Compatível

Tipo de tabela	Operation	Comando SQL Write	Status
MAIS	ALTER	ALTERAR TABELA - ALTERAR COLUNA	Compatível
MAIS	ALTER	ALTERAR TABELA - ADICIONAR COLUNAS	Compatível
MAIS	ALTER	ALTERAR TABELA - ADICIONAR PARTIÇÃO	Compatível
MAIS	ALTER	ALTERAR TABELA - SOLTAR PARTIÇÃO	Compatível
MAIS	ALTER	ALTERAR TABELA - RECUPERAR PARTIÇÕES	Compatível
MAIS	ALTER	REPARAR PARTIÇÕES DE SINCRONIZAÇÃO DE TABELAS	Compatível
DROP	DROP	DESCARTAR TABELA	Compatível

Tipo de tabela	Operation	Comando SQL Write	Status
MAIS	DROP	TABELA SUSPENSA - PURGA	Compatível
CREATE	CREATE	CRIAR TABELA - Gerenciado	Compatível
MAIS	CREATE	CRIAR TABELA - PARTICIONAR POR	Compatível
MAIS	CREATE	CRIAR TABELA SE NÃO EXISTIR	Compatível
MAIS	CREATE	CREATE TABLE LIKE	Compatível
MAIS	CREATE	CRIAR TABELA COMO SELEÇÃO	Compatível
CREATE	CREATE	CRIAR TABELA com LOCALIZAÇÃO - Tabela externa	Não suportado
QUADRO DE DADOS (UPSERT)	QUADRO DE DADOS (UPSERT)	saveAsTab le.Sobrescrever	Compatível
MAIS	QUADRO DE DADOS (UPSERT)	saveAsTab le.Anexar	Não suportado
MAIS	QUADRO DE DADOS (UPSERT)	saveAsTab le.Ignorar	Compatível

Tipo de tabela	Operation	Comando SQL Write	Status
MAIS	QUADRO DE DADOS (UPSERT)	saveAsTable. ErrorIfExists	Compatível
MAIS	QUADRO DE DADOS (UPSERT)	saveAsTable - Tabela externa (Caminho)	Não suportado
MAIS	QUADRO DE DADOS (UPSERT)	salvar (caminho) - DF v1	Não suportado
DATAFRAME (EXCLUIR)	DATAFRAME (EXCLUIR)	saveAsTable. Anexar	Não suportado
MAIS	DATAFRAME (EXCLUIR)	saveAsTable - Tabela externa (Caminho)	Não suportado
MAIS	DATAFRAME (EXCLUIR)	salvar (caminho) - DF v1	Não suportado
QUADRO DE DADOS (BULK_INSERT)	QUADRO DE DADOS (BULK_INSERT)	saveAsTable. Sobrescrever	Compatível
MAIS	QUADRO DE DADOS (BULK_INSERT)	saveAsTable. Anexar	Não suportado
MAIS	QUADRO DE DADOS (BULK_INSERT)	saveAsTable. Ignorar	Compatível
MAIS	QUADRO DE DADOS (BULK_INSERT)	saveAsTable. ErrorIfExists	Compatível
MAIS	QUADRO DE DADOS (BULK_INSERT)	saveAsTable - Tabela externa (Caminho)	Não suportado

Tipo de tabela	Operation	Comando SQL Write	Status
MAIS	QUADRO DE DADOS (BULK_INSERT)	salvar (caminho) - DF v1	Não suportado

Usando o EMR Serverless com para controle de acesso refinado AWS Lake Formation

Visão geral do

Com as versões 7.2.0 e superiores do Amazon EMR, aproveite AWS Lake Formation para aplicar controles de acesso refinados em tabelas do catálogo de dados que são apoiadas pelo S3. Esse recurso permite configurar controles de acesso em nível de tabela, linha, coluna e célula para consultas de read nos trabalhos do Spark no Amazon EMR Sem Servidor. Para configurar um controle de acesso refinado para trabalhos em lote e sessões interativas do Apache Spark, use o EMR Studio. Consulte as seções a seguir para saber mais sobre o Lake Formation e como usá-lo com o EMR Sem Servidor.

O uso do Amazon EMR Serverless incorre em cobranças adicionais. AWS Lake Formation Para obter mais informações, consulte [Preço do Amazon EMR](#).

Como o EMR Serverless funciona com AWS Lake Formation

Usar o EMR Sem Servidor com o Lake Formation permite impor uma camada de permissões em cada trabalho do Spark para aplicar o controle de permissões do Lake Formation quando o EMR Sem Servidor executa trabalhos. O EMR Sem Servidor usa [perfis de recursos do Spark](#) para criar dois perfis que executam trabalhos com eficiência. O perfil do usuário executa o código fornecido pelo usuário, enquanto o perfil do sistema impõe as políticas do Lake Formation. Para obter mais informações, consulte [O que é AWS Lake Formation](#) e [Considerações e limitações](#).

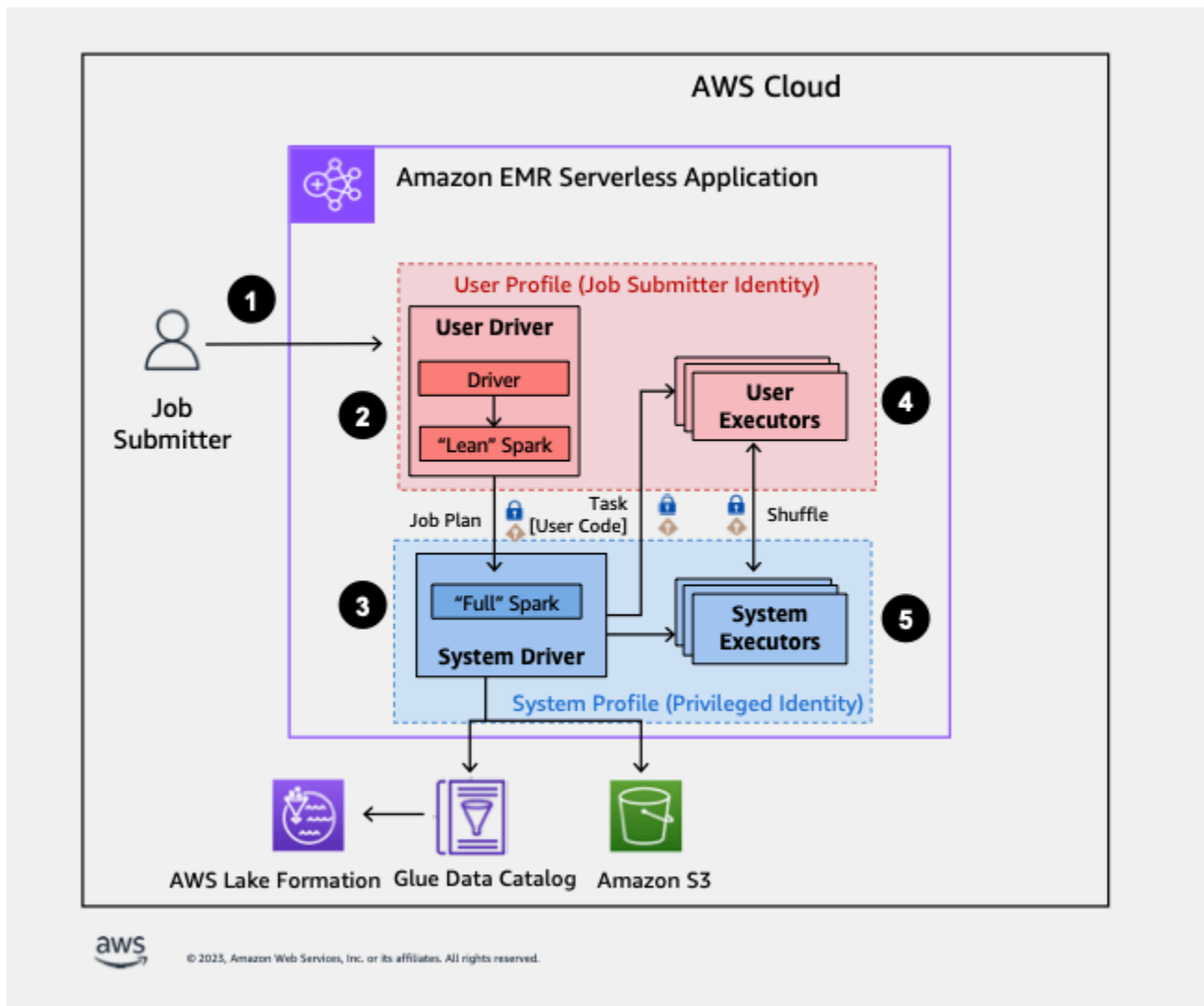
Ao usar a capacidade pré-inicializada com o Lake Formation, sugerimos que você tenha no mínimo dois drivers do Spark. Cada trabalho habilitado para o Lake Formation utiliza dois drivers do Spark, um para o perfil do usuário e outro para o perfil do sistema. Para obter a melhor performance, use o dobro do número de drivers para trabalhos habilitados para o Lake Formation em comparação com aqueles que não usam o Lake Formation.

Ao executar trabalhos do Spark no EMR Sem Servidor, considere também o impacto da alocação dinâmica no gerenciamento de recursos e na performance do cluster. A configuração `spark.dynamicAllocation.maxExecutors` do número máximo de executores por perfil de recurso se aplica aos executores do usuário e do sistema. Se você configurar esse número para ser igual ao número máximo permitido de executores, a execução do trabalho poderá ficar paralisada devido a um tipo de executor que usa todos os recursos disponíveis, o que impede o outro executor ao executar trabalhos.

Para que você não fique sem recursos, o EMR Sem Servidor define o número máximo padrão de executores por perfil de recurso como 90% do valor de `spark.dynamicAllocation.maxExecutors`. Você pode substituir essa configuração ao especificar `spark.dynamicAllocation.maxExecutorsRatio` com um valor entre 0 e 1. Além disso, configure as seguintes propriedades para otimizar a alocação de recursos e a performance geral:

- `spark.dynamicAllocation.cachedExecutorIdleTimeout`
- `spark.dynamicAllocation.shuffleTracking.timeout`
- `spark.cleaner.periodicGC.interval`

Confira a seguir uma visão geral de alto nível sobre como o EMR Sem Servidor obtém acesso aos dados protegidos pelas políticas de segurança do Lake Formation.



1. Um usuário envia uma tarefa do Spark para um aplicativo AWS Lake Formation EMR Serverless habilitado.
2. O EMR Sem Servidor envia o trabalho para um driver de usuário e executa o trabalho no perfil do usuário. O driver do usuário executa uma versão enxuta do Spark que não tem a capacidade de iniciar tarefas, solicitar executores, acessar o S3 ou o Glue Catalog. Ele cria um plano de trabalho.
3. O EMR Sem Servidor configura um segundo driver chamado driver do sistema e o executa no perfil do sistema (com uma identidade privilegiada). O EMR Sem Servidor configura um canal TLS criptografado entre os dois drivers para comunicação. O driver do usuário usa o canal para enviar os planos de trabalho ao driver do sistema. O driver do sistema não executa o código enviado pelo usuário. Ele executa o Spark completo e se comunica com o S3 e com o Data

Catalog para acesso aos dados. Ele solicita executores e compila o plano de trabalho em uma sequência de estágios de execução.

4. Em seguida, o EMR Sem Servidor executa os estágios nos executores com o driver do usuário ou o driver do sistema. O código do usuário em qualquer estágio é executado exclusivamente nos executores do perfil do usuário.
5. Os estágios que lêem dados das tabelas do Catálogo de Dados protegidas por AWS Lake Formation ou que aplicam filtros de segurança são delegados aos executores do sistema.

Como habilitar o Lake Formation no Amazon EMR

Para habilitar o Lake Formation, defina `spark.emr-serverless.lakeformation.enabled` como `true` na classificação `spark-defaults` do parâmetro de configuração de runtime ao [criar uma aplicação do EMR Sem Servidor](#).

```
aws emr-serverless create-application \  
  --release-label emr-7.12.0 \  
  --runtime-configuration '{  
    "classification": "spark-defaults",  
    "properties": {  
      "spark.emr-serverless.lakeformation.enabled": "true"  
    }  
  }' \  
  --type "SPARK"
```

Você também pode habilitar o Lake Formation ao criar uma aplicação no EMR Studio. Escolha Usar Lake Formation para um controle de acesso refinado, disponível em Configurações adicionais.

A [criptografia entre trabalhadores](#) é habilitada por padrão ao usar o Lake Formation com o EMR Sem Servidor, então você não precisa habilitar explicitamente a criptografia entre trabalhadores novamente.

Como habilitar o Lake Formation para trabalhos no Spark

Para habilitar o Lake Formation em trabalhos individuais do Spark, defina `spark.emr-serverless.lakeformation.enabled` como `verdadeiro` ao usar `spark-submit`.

```
--conf spark.emr-serverless.lakeformation.enabled=true
```

Permissões do IAM do perfil de runtime do trabalho

As permissões do Lake Formation controlam o acesso aos recursos do AWS Glue Data Catalog, aos locais do Amazon S3 e aos dados subjacentes nesses locais. As permissões do IAM controlam o acesso ao Lake Formation, ao AWS Glue APIs e aos recursos. Embora você possa ter a permissão do Lake Formation para acessar uma tabela no Data Catalog (SELECT), a operação falhará se você não tiver a permissão do IAM na operação de API `glue:Get*`.

Confira a seguir um exemplo de política de como fornecer permissões do IAM para acesso a um script no S3, upload de logs no S3, permissões da API do AWS Glue e permissão para acessar o Lake Formation.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "ScriptAccess",
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3::*.amzn-s3-demo-bucket/scripts",
        "arn:aws:s3::*.amzn-s3-demo-bucket/*"
      ]
    },
    {
      "Sid": "LoggingAccess",
      "Effect": "Allow",
      "Action": [
        "s3:PutObject"
      ],
      "Resource": [
        "arn:aws:s3::amzn-s3-demo-bucket/logs/*"
      ]
    },
    {
      "Sid": "GlueCatalogAccess",
      "Effect": "Allow",
```

```
"Action": [
  "glue:Get*",
  "glue:Create*",
  "glue:Update*"
],
"Resource": [
  "*"
]
},
{
  "Sid": "LakeFormationAccess",
  "Effect": "Allow",
  "Action": [
    "lakeformation:GetDataAccess"
  ],
  "Resource": [
    "*"
  ]
}
]
```

Configuração de permissões do Lake Formation para perfil de runtime do trabalho

Primeiro, registre a localização da tabela do Hive no Lake Formation. Em seguida, crie permissões para o perfil de runtime do trabalho na tabela desejada. Para obter mais detalhes sobre Lake Formation, consulte [O que é AWS Lake Formation?](#) no Guia do AWS Lake Formation desenvolvedor.

Depois de configurar as permissões do Lake Formation, envie trabalhos do Spark no Amazon EMR Sem Servidor. Para obter mais informações sobre os trabalhos do Spark, consulte os [Exemplos do Spark](#).

Envio da execução de um trabalho

Depois de concluir a configuração das concessões do Lake Formation, você pode [enviar trabalhos do Spark no EMR Sem Servidor](#). A seção a seguir mostra exemplos de como configurar e enviar propriedades de execução de trabalho.

Requisitos de permissão

Tabelas não registradas em AWS Lake Formation

Para tabelas não registradas AWS Lake Formation, a função de tempo de execução do trabalho acessa o AWS Glue Data Catalog e os dados da tabela subjacente no Amazon S3. Isso exige que a função de tempo de execução do trabalho tenha as permissões apropriadas do IAM para as operações do AWS Glue e do Amazon S3.

Tabelas registradas em AWS Lake Formation

Para tabelas registradas com AWS Lake Formation, a função de tempo de execução do trabalho acessa os metadados do AWS Glue Data Catalog, enquanto as credenciais temporárias fornecidas pela Lake Formation acessam os dados da tabela subjacente no Amazon S3. As permissões do Lake Formation necessárias para executar uma operação dependem do AWS Glue Data Catalog e das chamadas de API do Amazon S3 iniciadas pelo trabalho do Spark e podem ser resumidas da seguinte forma:

- A permissão DESCRIBE permite que a função de tempo de execução leia os metadados da tabela ou do banco de dados no Catálogo de Dados
- A permissão ALTER permite que a função de tempo de execução modifique os metadados da tabela ou do banco de dados no Catálogo de Dados
- A permissão DROP permite que a função de tempo de execução exclua metadados da tabela ou do banco de dados do Catálogo de Dados
- A permissão SELECT permite que a função de tempo de execução leia dados da tabela do Amazon S3
- A permissão INSERT permite que a função de tempo de execução grave dados de tabela no Amazon S3
- A permissão DELETE permite que a função de tempo de execução exclua dados da tabela do Amazon S3

Note

O Lake Formation avalia as permissões preguiçosamente quando um trabalho do Spark chama o AWS Glue para recuperar os metadados da tabela e o Amazon S3 para recuperar os dados da tabela. Trabalhos que usam uma função de tempo de execução com

permissões insuficientes não falharão até que o Spark faça uma chamada ao AWS Glue ou ao Amazon S3 que exija a permissão ausente.

Note

Na seguinte matriz de tabela suportada:

- As operações marcadas como Suportadas usam exclusivamente as credenciais do Lake Formation para acessar os dados das tabelas registradas no Lake Formation. Se as permissões do Lake Formation forem insuficientes, a operação não retornará às credenciais da função de tempo de execução. Para tabelas não registradas no Lake Formation, as credenciais da função de tempo de execução do trabalho acessam os dados da tabela.
- As operações marcadas como compatíveis com permissões do IAM na localização do Amazon S3 não usam as credenciais do Lake Formation para acessar os dados da tabela subjacente no Amazon S3. Para executar essas operações, a função de tempo de execução do trabalho deve ter as permissões necessárias do Amazon S3 IAM para acessar os dados da tabela, independentemente de a tabela estar registrada no Lake Formation.

Hive

Operation	AWS Lake Formation permissões	Status do suporte
SELECT	SELECT	Compatível
CRIAR TABELA	CRIAR_TABELA	Compatível
CREATE TABLE LIKE	CRIAR_TABELA	Compatível com permissões do IAM na localização do Amazon S3
CRIAR TABELA COMO SELEÇÃO	CRIAR_TABELA	Compatível com permissões do IAM na localização do Amazon S3

Operation	AWS Lake Formation permissões	Status do suporte
DESCRIBE TABLE	DESCRIBE	Compatível
SHOW TBLPROPERTIES	DESCRIBE	Compatível
SHOW COLUMNS	DESCRIBE	Compatível
SHOW PARTITIONS	DESCRIBE	Compatível
SHOW CREATE TABLE	DESCRIBE	Compatível
ALTERAR TABELA tablename	SELECIONAR e ALTERAR	Compatível
ALTERAR LOCALIZAÇÃO DO tablename CONJUNTO DE TABELAS	-	Não compatível
ALTERAR TABELA tablename ADICIONAR PARTIÇÃO	SELECIONAR, INSERIR e ALTERAR	Compatível
REPAIR TABLE	SELECIONAR e ALTERAR	Compatível
CARREGAR DADOS		Não compatível
INSERT	INSERIR e ALTERAR	Compatível
INSERT OVERWRITE	SELECIONAR, INSERIR, EXCLUIR e ALTERAR	Compatível
DESCARTAR TABELA	SELECIONAR, DESCARTAR, EXCLUIR e ALTERAR	Compatível

Operation	AWS Lake Formation permissões	Status do suporte
TRUNCATE TABLE	SELECIONAR, INSERIR, EXCLUIR e ALTERAR	Compatível
Gravador de quadros de dados V1	Igual à operação SQL correspondente	Compatível ao anexar dados a uma tabela existente. Consulte as considerações e limitações para obter mais informações.
Gravador de quadros de dados V2	Igual à operação SQL correspondente	Compatível ao anexar dados a uma tabela existente. Consulte as considerações e limitações para obter mais informações.

Iceberg

Operation	AWS Lake Formation permissões	Status do suporte
SELECT	SELECT	Compatível
CRIAR TABELA	CRIAR_TABELA	Compatível
CREATE TABLE LIKE	CRIAR_TABELA	Compatível com permissões do IAM na localização do Amazon S3
CRIAR TABELA COMO SELEÇÃO	CRIAR_TABELA	Compatível com permissões do IAM na localização do Amazon S3
SUBSTITUA A TABELA CONFORME SELECIONADO	SELECIONAR, INSERIR e ALTERAR	Compatível
DESCRIBE TABLE	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3

Operation	AWS Lake Formation permissões	Status do suporte
SHOW TBLPROPERTIES	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3
SHOW CREATE TABLE	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3
ALTER TABLE	SELECIONAR, INSERIR e ALTERAR	Compatível
ALTER TABLE SET LOCATION	SELECIONAR, INSERIR e ALTERAR	Compatível com permissões do IAM na localização do Amazon S3
ALTERAR A GRAVAÇÃO DA TABELA ORDENADA POR	SELECIONAR, INSERIR e ALTERAR	Compatível com permissões do IAM na localização do Amazon S3
ALTERE A GRAVAÇÃO DA TABELA DISTRIBUÍDA POR	SELECIONAR, INSERIR e ALTERAR	Compatível com permissões do IAM na localização do Amazon S3
ALTERAR TABELA RENOMEAR TABELA	CREATE_TABLE e DROP	Compatível
INSERT INTO	SELECIONAR, INSERIR e ALTERAR	Compatível
INSERT OVERWRITE	SELECIONAR, INSERIR e ALTERAR	Compatível
DELETE	SELECIONAR, INSERIR e ALTERAR	Compatível
UPDATE	SELECIONAR, INSERIR e ALTERAR	Compatível

Operation	AWS Lake Formation permissões	Status do suporte
MERGE INTO	SELECIONAR, INSERIR e ALTERAR	Compatível
DESCARTAR TABELA	SELECIONAR, EXCLUIR e SOLTAR	Compatível
DataFrame Escritor V1	-	Não compatível
DataFrame Escritor V2	Igual à operação SQL correspondente	Compatível ao anexar dados a uma tabela existente. Consulte as considerações e limitações para obter mais informações.
Tabelas de metadados	SELECT	Compatível. Certas tabelas estão ocultas. Consulte as considerações e limitações para obter mais informações.
Procedimentos armazenados	-	<p>Compatível com tabelas que atendam às seguintes condições:</p> <ul style="list-style-type: none"> • Tabelas não registradas em AWS Lake Formation • Tabelas que não usam <code>register_table</code> e <code>migrate</code> <p>Consulte as considerações e limitações para obter mais informações.</p>

Configuração do Spark para o Iceberg: o exemplo a seguir mostra como configurar o Spark com o Iceberg. Para executar trabalhos do Iceberg, é necessário fornecer as propriedades `spark-submit` a seguir.

```
--conf spark.sql.catalog.spark_catalog=org.apache.iceberg.spark.SparkSessionCatalog
```

```
--conf spark.sql.catalog.spark_catalog.warehouse=<S3_DATA_LOCATION>
--conf spark.sql.catalog.spark_catalog.glue.account-id=<ACCOUNT_ID>
--conf spark.sql.catalog.spark_catalog.client.region=<REGION>
--conf spark.sql.catalog.spark_catalog.glue.endpoint=https://
glue.<REGION>.amazonaws.com
```

Hudi

Operation	AWS Lake Formation permissões	Status do suporte
SELECT	SELECT	Compatível
CRIAR TABELA	CRIAR_TABELA	Compatível com permissões do IAM na localização do Amazon S3
CREATE TABLE LIKE	CRIAR_TABELA	Compatível com permissões do IAM na localização do Amazon S3
CRIAR TABELA COMO SELEÇÃO	-	Não compatível
DESCRIBE TABLE	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3
SHOW TBLPROPERTIES	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3
SHOW COLUMNS	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3
SHOW CREATE TABLE	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3

Operation	AWS Lake Formation permissões	Status do suporte
ALTER TABLE	SELECT	Compatível com permissões do IAM na localização do Amazon S3
INSERT INTO	SELECIONAR e ALTERAR	Compatível com permissões do IAM na localização do Amazon S3
INSERT OVERWRITE	SELECIONAR e ALTERAR	Compatível com permissões do IAM na localização do Amazon S3
DELETE	-	Não compatível
UPDATE	-	Não compatível
MERGE INTO	-	Não compatível
DESCARTAR TABELA	SELECIONAR e SOLTAR	Compatível com permissões do IAM na localização do Amazon S3
DataFrame Escritor V1	-	Não compatível
DataFrame Escritor V2	Igual à operação SQL correspondente	Compatível com permissões do IAM na localização do Amazon S3
Tabelas de metadados	-	Não compatível
Manutenção de tabelas e recursos utilitários	-	Não compatível

As amostras a seguir configuram o Spark com o Hudi, especificando os locais dos arquivos e outras propriedades necessárias para o uso.

Configuração do Spark para Hudi: quando usado em um notebook, esse trecho especifica o caminho para o arquivo JAR do pacote Hudi Spark, que habilita a funcionalidade Hudi no Spark. Ele também configura o Spark para usar o AWS Glue Data Catalog como metastore.

```
%%configure -f
{
  "conf": {
    "spark.jars": "/usr/lib/hudi/hudi-spark-bundle.jar",
    "spark.hadoop.hive.metastore.client.factory.class":
"com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalogHiveClientFactory",
    "spark.serializer": "org.apache.spark.serializer.JavaSerializer",
    "spark.sql.catalog.spark_catalog":
"org.apache.spark.sql.hudi.catalog.HoodieCatalog",
    "spark.sql.extensions":
"org.apache.spark.sql.hudi.HoodieSparkSessionExtension"
  }
}
```

Configuração do Spark para Hudi with AWS Glue: esse trecho, quando usado em um notebook, habilita o Hudi como um formato de lago de dados compatível e garante que as bibliotecas e dependências do Hudi estejam disponíveis.

```
%%configure
{
  "--conf": "spark.serializer=org.apache.spark.serializer.JavaSerializer --conf
spark.sql.catalog.spark_catalog=org.apache.spark.sql.hudi.catalog.HoodieCatalog --
conf
spark.sql.extensions=org.apache.spark.sql.hudi.HoodieSparkSessionExtension",
  "--datalake-formats": "hudi",
  "--enable-glue-datacatalog": True,
  "--enable-lakeformation-fine-grained-access": "true"
}
```

Delta Lake

Operation	AWS Lake Formation permissões	Status do suporte
SELECT	SELECT	Compatível
CRIAR TABELA	CRIAR_TABELA	Compatível

Operation	AWS Lake Formation permissões	Status do suporte
CREATE TABLE LIKE	-	Não compatível
CRIAR TABELA COMO SELEÇÃO	CRIAR_TABELA	Compatível
SUBSTITUA A TABELA CONFORME SELECIONADO	SELECIONAR, INSERIR e ALTERAR	Compatível
DESCRIBE TABLE	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3
SHOW TBLPROPERTIES	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3
SHOW COLUMNS	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3
SHOW CREATE TABLE	DESCRIBE	Compatível com permissões do IAM na localização do Amazon S3
ALTER TABLE	SELECIONAR e INSERIR	Compatível
ALTER TABLE SET LOCATION	SELECIONAR e INSERIR	Compatível com permissões do IAM na localização do Amazon S3
ALTERAR tablename AGRUPAMENTO DE TABELAS POR	SELECIONAR e INSERIR	Compatível com permissões do IAM na localização do Amazon S3


```
--conf
spark.sql.catalog.spark_catalog=org.apache.spark.sql.delta.catalog.DeltaCatalog \
```

Depurar trabalhos

Note

Com esse recurso, acesse os logs stdout e stderr dos trabalhadores do perfil do sistema que podem conter informações confidenciais e não filtradas. A permissão a seguir deve ser usada somente para acessar dados que não sejam de produção. Para aplicações criadas para uso com trabalhos de produção, sugerimos que você adicione essas permissões apenas a administradores ou usuários com acesso elevado aos dados.

Com o EMR-7.3.0 e versões posteriores, o EMR Sem Servidor está habilitando a capacidade de autodepuração para trabalhos em lotes habilitados para Lake Formation. Para fazer isso, use o novo parâmetro `accessSystemProfileLogs` na [GetDashboardForJobRunAPI](#). Se `accessSystemProfileLogs` estiver definido como verdadeiro, você poderá acessar os registros stdout e stderr dos trabalhadores do perfil do sistema, que podem ser usados para depurar um trabalho em lote do EMR Serverless habilitado para Lake Formation.

```
aws emr-serverless get-dashboard-for-job-run \
  --application-id application-id
  --job-run-id job-run-id
  --access-system-profile-logs
```

Permissões obrigatórias

O diretor que deseja depurar trabalhos em lote habilitados para Lake Formation usando `GetDashboardForJobRun` deve ter as seguintes permissões adicionais:

```
{
  "Sid": "AccessSystemProfileLogs",
  "Effect": "Allow",
  "Action": [
    "emr-serverless:GetDashboardForJobRun",
    "emr-serverless:AccessSystemProfileLogs",
    "glue:GetDatabases",
```

```
    "glue:SearchTables"
  ],
  "Resource": [
    "arn:aws:emr-serverless:region:account-id:/applications/applicationId/
jobruns/jobid",
    "arn:aws:glue:region:account-id:catalog",
    "arn:aws:glue:region:account-id:database/*",
    "arn:aws:glue:region:account-id:table/*/*"
  ]
}
```

Considerações

Os logs de perfil do sistema para depuração são visíveis para trabalhos que acessam bancos de dados ou tabelas no Lake Formation na mesma conta do trabalho. Eles não são visíveis nos seguintes cenários:

- Se o catálogo de dados gerenciado usando as permissões do Lake Formation tiver bancos de dados e tabelas entre contas
- Se o catálogo de dados gerenciado usando as permissões do Lake Formation tiver links de recursos

Trabalhar com visualizações do Catálogo de dados do Glue

Você pode criar e gerenciar visualizações no AWS Glue Data Catalog para uso com o EMR Serverless. Elas são comumente conhecidas como visualizações do AWS Glue Data Catalog. Essas visualizações são úteis porque oferecem suporte a vários mecanismos de consulta SQL, para que você possa acessar a mesma visualização em diferentes AWS serviços, como EMR Serverless e Amazon Athena Amazon Redshift.

Ao criar uma exibição no Catálogo de Dados, use concessões de recursos e controles de acesso baseados em tags AWS Lake Formation para conceder acesso a ela. Usando esse método de controle de acesso, você não precisa configurar acesso adicional às tabelas referenciadas ao criar a exibição. Esse método de concessão de permissões é chamado de semântica definidora, e essas exibições são chamadas de visualizações definidoras. Para obter mais informações sobre o controle de acesso no Lake Formation, consulte [Conceder e revogar permissões nos recursos do Catálogo de Dados no](#) Guia do Desenvolvedor do AWS Lake Formation.

As visualizações do Data Catalog são úteis para os seguintes casos de uso:

- **Controle de acesso granular:** é possível criar uma visualização que restringe o acesso aos dados com base nas permissões necessárias ao usuário. Por exemplo, você pode usar as exibições do Data Catalog para evitar que funcionários que não trabalham no departamento de RH vejam informações de identificação pessoal (PII).
- **Definição completa da exibição:** ao aplicar filtros à sua visualização no Catálogo de dados, você garante que os registros de dados em uma exibição do Catálogo de dados estejam sempre completos.
- **Segurança aprimorada:** a definição da consulta usada para criar a exibição deve estar completa. Esse benefício significa que as visualizações no Data Catalog não são suscetíveis a comandos SQL de acessos mal-intencionados.
- **Compartilhamento simples de dados** — compartilhe dados com outras AWS contas sem mover dados. Para obter mais informações, consulte [Compartilhamento de dados entre contas no Lake Formation](#).

Criação de uma visualização do Catálogo de Dados

Há maneiras diferentes de criar uma exibição do Catálogo de dados. Isso inclui o uso do AWS CLI ou do Spark SQL. Veja alguns exemplos.

Using SQL

Veja a seguir a sintaxe para a criação de uma exibição do Catálogo de dados. Observe o tipo de visualização `MULTI DIALECT`. Isso distingue a exibição do Catálogo de dados de outras visualizações. O predicado `SECURITY` é especificado como `DEFINER`. Isso indica uma visualização do Catálogo de dados com semântica `DEFINER`.

```
CREATE [ OR REPLACE ] PROTECTED MULTI DIALECT VIEW [IF NOT EXISTS] view_name
[(column_name [COMMENT column_comment], ...) ]
[ COMMENT view_comment ]
[TBLPROPERTIES (property_name = property_value, ... )]
SECURITY DEFINER
AS query;
```

Veja a seguir um exemplo de instrução `CREATE`, seguindo a sintaxe:

```
CREATE PROTECTED MULTI DIALECT VIEW catalog_view
SECURITY DEFINER
AS
```

```
SELECT order_date, sum(totalprice) AS price
FROM source_table
GROUP BY order_date
```

Você também pode criar uma exibição no modo de simulação, usando SQL, para testar a criação da exibição, sem realmente criar o recurso. O uso dessa opção resulta em uma “execução seca” que valida a entrada e, se a validação for bem-sucedida, retorna o JSON do objeto da tabela AWS Glue que representará a visualização. Nesse caso, a visualização real não é criada.

```
CREATE [ OR REPLACE ] PROTECTED MULTI DIALECT VIEW view_name
SECURITY DEFINER
[ SHOW VIEW JSON ]
AS view-sql
```

Using the AWS CLI

Note

Quando você usa o comando da CLI, o SQL usado para criar a exibição não é analisado. Isso pode resultar em um caso em que a exibição é criada, mas as consultas não são bem-sucedidas. Certifique-se de testar sua sintaxe SQL antes de criar a exibição.

Use o seguinte comando da CLI para criar uma visualização:

```
aws glue create-table --cli-input-json '{
  "DatabaseName": "database",
  "TableInput": {
    "Name": "view",
    "StorageDescriptor": {
      "Columns": [
        {
          "Name": "col1",
          "Type": "data-type"
        },
        ...
        {
          "Name": "col_n",
          "Type": "data-type"
        }
      ]
    }
  }
},
```

```

    "SerdeInfo": {}
  },
  "ViewDefinition": {
    "SubObjects": [
      "arn:aws:glue:aws-region:aws-account-id:table/database/referenced-table1",
      ...
      "arn:aws:glue:aws-region:aws-account-id:table/database/referenced-tableN",
    ],
    "IsProtected": true,
    "Representations": [
      {
        "Dialect": "SPARK",
        "DialectVersion": "1.0",
        "ViewOriginalText": "Spark-SQL",
        "ViewExpandedText": "Spark-SQL"
      }
    ]
  }
}
}'

```

Operações compatíveis para visualizações

Os fragmentos de comandos apresentados a seguir mostram como é possível trabalhar com as visualizações do Catálogo de Dados de diferentes formas:

- CREATE VIEW

Cria uma visualização do Catálogo de Dados. A seguir, um exemplo que demonstra a criação da visualização usando uma tabela existente:

```

CREATE PROTECTED MULTI DIALECT VIEW catalog_view
SECURITY DEFINER AS SELECT * FROM my_catalog.my_database.source_table

```

- ALTER VIEW

Sintaxe disponível:

- ALTER VIEW *view_name* [FORCE] ADD DIALECT AS *query*
- ALTER VIEW *view_name* [FORCE] UPDATE DIALECT AS *query*
- ALTER VIEW *view_name* DROP DIALECT

É possível usar a opção `FORCE ADD DIALECT` para aplicar a atualização do esquema e dos objetos secundários de acordo com o novo dialeto do mecanismo. Lembre-se de que essa ação pode resultar em erros de consulta, caso a opção `FORCE` não seja usada para atualizar os demais dialetos do mecanismo. A seguir, apresentamos um exemplo:

```
ALTER VIEW catalog_view FORCE ADD DIALECT
AS
SELECT order_date, sum(totalprice) AS price
FROM source_table
GROUP BY orderdate;
```

A seguir, é demonstrado como alterar uma visualização para atualizar o dialeto:

```
ALTER VIEW catalog_view UPDATE DIALECT AS
SELECT count(*) FROM my_catalog.my_database.source_table;
```

- **DESCRIBE VIEW**

Sintaxe disponível para descrever uma visualização:

- `SHOW COLUMNS {FROM|IN} view_name [{FROM|IN} database_name]`— Se o usuário tiver as permissões necessárias de AWS Glue e Lake Formation para descrever a visualização, ele poderá listar as colunas. A seguir, alguns exemplos de comandos que demonstram como visualizar colunas:

```
SHOW COLUMNS FROM my_database.source_table;
SHOW COLUMNS IN my_database.source_table;
```

- `DESCRIBE view_name`— Se o usuário tiver as permissões necessárias de AWS Glue e Lake Formation para descrever a visualização, ele poderá listar as colunas na exibição junto com seus metadados.

- **DROP VIEW**

Sintaxe disponível:

- `DROP VIEW [IF EXISTS] view_name`

A seguir, um exemplo de instrução `DROP` que demonstra como realizar um teste para verificar a existência da visualização antes de sua remoção:

```
DROP VIEW IF EXISTS catalog_view;
```

- **MOSTRAR CRIAR VISUALIZAÇÃO**

- `SHOW CREATE VIEW view_name`: mostra a instrução SQL responsável pela criação da visualização especificada. A seguir, um exemplo que demonstra o processo de criação de uma visualização no Catálogo de Dados:

```
SHOW CREATE TABLE my_database.catalog_view;  
CREATE PROTECTED MULTI DIALECT VIEW my_catalog.my_database.catalog_view (  
  net_profit,  
  customer_id,  
  item_id,  
  sold_date)  
TBLPROPERTIES (  
  'transient_lastDdlTime' = '1736267222')  
SECURITY DEFINER AS SELECT * FROM  
my_database.store_sales_partitioned_1f WHERE customer_id IN (SELECT customer_id  
from source_table limit 10)
```

- **SHOW VIEWS**

Lista todas as visualizações existentes no catálogo, incluindo as visualizações regulares, as visualizações com múltiplos dialetos (MDV, na sigla em inglês) e as MDVs sem suporte ao dialeto Spark. A sintaxe disponível é a seguinte:

- `SHOW VIEWS [{ FROM | IN } database_name] [LIKE regex_pattern]:`

A seguir, um exemplo que demonstra um comando para mostrar as visualizações:

```
SHOW VIEWS IN marketing_analytics LIKE 'catalog_view*';
```

Para obter mais informações sobre como criar e configurar visualizações do catálogo de dados, consulte Building [AWS Glue Data Catalog no Developer](#) Guide. AWS Lake Formation

Consulta de uma visualização do Catálogo de Dados

Depois de criar uma visualização do catálogo de dados, você pode consultá-la usando uma tarefa do Amazon EMR Serverless Spark que tenha AWS Lake Formation um controle de acesso refinado ativado. O perfil de runtime do trabalho deve ter a permissão `SELECT` do Lake Formation na

visualização do Catálogo de dados. Não é necessário conceder acesso às tabelas subjacentes referenciadas na visualização.

Depois de configurar tudo, você pode consultar sua exibição. Por exemplo, depois de criar uma aplicação do EMR Sem Servidor no EMR Studio, execute a consulta a seguir para acessar uma exibição.

```
SELECT * from my_database.catalog_view LIMIT 10;
```

Uma função útil é `invoker_principal`. Ela retorna o identificador exclusivo do perfil de runtime do trabalho do EMRS. Isso pode ser usado para controlar a saída da visualização, com base na entidade principal de invocação. Você pode usar isso para adicionar uma condição à sua visualização que refina os resultados da consulta, com base no perfil de chamada. O perfil de runtime do trabalho deve ter permissão para a ação `LakeFormation:GetDataLakePrincipal` do IAM usar essa função.

```
select invoker_principal();
```

Você pode adicionar essa função a uma cláusula `WHERE`, por exemplo, para refinar os resultados da consulta.

Considerações e limitações

Quando você cria exibições do Catálogo de dados, o seguinte se aplica:

- Você só pode criar exibições do Catálogo de dados com o Amazon EMR 7.6 e versões superiores.
- O responsável pela definição da visualização do Catálogo de Dados deve ter permissão de acesso `SELECT` nas tabelas de base subjacentes acessadas pela visualização. A operação de criação da visualização do Catálogo de Dados não será bem-sucedida caso uma das tabelas de base tenha filtros do Lake Formation atribuídos ao perfil “definer”.
- As tabelas de base não devem ter a permissão de data lake `IAMAllowedPrincipals` no Lake Formation. Se presente, o erro de visualizações de vários dialetos só pode fazer referência a tabelas sem a permissão do `IAMAllowed Principal`.
- A localização da tabela no Amazon S3 deve ser registrada como uma localização de data lake do Lake Formation. Se a tabela não estiver registrada, o erro Visualizações de vários dialetos só podem fazer referência a tabelas gerenciadas pelo Lake Formation ocorrerá. Para obter informações sobre como registrar locais do Amazon S3 em Lake Formation, consulte [Registro de um local do Amazon S3](#) no Guia do desenvolvedor. AWS Lake Formation

- Você só pode criar visualizações PROTECTED do Data Catalog. Não há suporte para exibições UNPROTECTED.
- Você não pode referenciar tabelas em outra AWS conta em uma definição de exibição do Catálogo de Dados. Além disso, não é possível referenciar uma tabela da mesma conta que esteja em uma região diferente.
- Para compartilhar dados entre contas ou regiões, a visualização inteira deve ser compartilhada entre contas e entre regiões, usando links de recurso do Lake Formation.
- As funções definidas pelo usuário (UDFs) não são suportadas.
- É possível usar visualizações baseadas em tabelas do Iceberg. Os formatos de tabela aberta Apache Hudi e Delta Lake também são compatíveis.
- Não é possível fazer referência a outras visualizações nas exibições do Data Catalog.
- Um esquema de visualização do AWS Glue Data Catalog é sempre armazenado em letras minúsculas. Por exemplo, se você usar uma instrução DDL para criar uma exibição do Catálogo de dados do Glue com uma coluna chamada Castle, a coluna criada no Catálogo de dados do Glue será transformada em minúsculas, para castle. Se você especificar o nome da coluna em uma consulta DML como Castle ou CASTLE, o EMR Spark mudará o nome para letras minúsculas para você executar a consulta. Porém, o título da coluna é exibido usando as letras maiúsculas e minúsculas que você especificou na consulta.

Se quiser que uma consulta falhe em um caso em que o nome da coluna especificado na consulta DML não corresponda ao nome da coluna no Catálogo de dados do Glue, defina `spark.sql.caseSensitive=true`.

Suporte ao formato de tabela aberta

O EMR Serverless suporta consultas SELECT no Apache Hive, Apache Iceberg, Delta Lake (7.6.0+) e Apache Hudi (7.6.0+). A partir do EMR 7.12, as operações DML e DDL que modificam os dados da tabela são suportadas nas tabelas Apache Hive, Apache Iceberg e Delta Lake usando credenciais vendidas do Lake Formation.

Considerações e limitações

Geral

Analise as seguintes limitações ao usar o Lake Formation com o EMR Serverless.

Note

Quando você habilita o Lake Formation para um trabalho do Spark no EMR Sem Servidor, o trabalho inicia um driver de sistema e um driver de usuário. Se você especificou a capacidade pré-inicializada na inicialização, os drivers são provisionados na capacidade pré-inicializada e o número de drivers do sistema é igual ao número de drivers de usuário que você especifica. Se você escolher a capacidade sob demanda, o EMR Sem Servidor iniciará um driver de sistema além de um driver de usuário. Para estimar os custos associados ao trabalho do EMR Sem Servidor com o Lake Formation, use o serviço [AWS Calculadora de Preços](#).

- O Amazon EMR Sem Servidor com Lake Formation está disponível em todas as [regiões do EMR Sem Servidor](#) com suporte.
- As aplicações habilitadas para Lake Formation não oferecem suporte ao uso de [imagens personalizadas do EMR Sem Servidor](#).
- Você não pode desativar `DynamicResourceAllocation` para trabalhos do Lake Formation.
- Você só pode usar o Lake Formation com trabalhos do Spark.
- O EMR Sem Servidor com Lake Formation oferece suporte apenas a uma única sessão do Spark durante todo o trabalho.
- O EMR Sem Servidor com Lake Formation só oferece suporte a consultas de tabelas entre contas compartilhadas por meio de links de recursos.
- As seguintes opções não são compatíveis:
 - Conjuntos de dados distribuídos resilientes (RDD)
 - Streaming do Spark
 - Controle de acesso para colunas aninhadas
- O EMR Sem Servidor bloqueia funcionalidades que podem prejudicar o isolamento completo do driver do sistema, incluindo as seguintes:
 - UDTs, Hive UDFs e qualquer função definida pelo usuário que envolva classes personalizadas
 - Fontes de dados personalizadas
 - Fornecimento de JARs adicionais para extensão, conector ou metastore do Spark
 - Comando `ANALYZE TABLE`

- [Se seu aplicativo EMR Serverless estiver em uma sub-rede privada com endpoints VPC para o Amazon S3 e você anexar uma política de endpoint para controlar o acesso, antes que seus trabalhos possam enviar dados de log para o Amazon S3 gerenciado, inclua as permissões detalhadas em Armazenamento gerenciado AWS em sua política de VPC para o endpoint do gateway S3.](#) Para solicitações de solução de problemas, entre em contato com AWS o suporte.
- A partir do Amazon EMR 7.9.0, o Spark FGAC oferece suporte ao AFile sistema S3 quando usado com o esquema s3a://.
- O Amazon EMR 7.11 oferece suporte à criação de tabelas gerenciadas usando CTAS.
- O Amazon EMR 7.12 oferece suporte à criação de tabelas gerenciadas e externas usando CTAS.

Permissões

- Para aplicar controles de acesso, as operações EXPLAIN PLAN e DDL, como DESCRIBE TABLE, não expõem informações restritas.
- Quando você registra a localização de uma tabela no Lake Formation, o acesso aos dados usa as credenciais armazenadas do Lake Formation em vez das permissões IAM da função de tempo de execução de trabalho do EMR Serverless. Os trabalhos falharão se a função registrada para a localização da tabela estiver configurada incorretamente, mesmo quando a função de tempo de execução tiver permissões do S3 IAM para essa localização.
- A partir do Amazon EMR 7.12, você pode gravar nas tabelas existentes do Hive e do Iceberg usando DataFrameWriter (V2) com as credenciais do Lake Formation no modo de acréscimo. Para operações de substituição ou ao criar novas tabelas, o EMR usa as credenciais da função de tempo de execução para modificar os dados da tabela.
- As seguintes limitações se aplicam ao usar exibições ou tabelas em cache como dados de origem (essas limitações não se aplicam às visualizações do AWS Glue Data Catalog):
 - Para operações MERGE, DELETE e UPDATE
 - Compatível: uso de visualizações e tabelas em cache como tabelas de origem.
 - Não suportado: uso de visualizações e tabelas em cache em cláusulas de atribuição e condição.
 - Para as operações CREATE OR REPLACE e REPLACE TABLE AS SELECT:
 - Não suportado: uso de visualizações e tabelas em cache como tabelas de origem.
- As tabelas Delta Lake com dados UDFs na fonte suportam as operações MERGE, DELETE e UPDATE somente quando o vetor de exclusão está ativado.

Registros e depuração

- O EMR Sem Servidor restringe o acesso aos logs do Spark do driver do sistema em aplicações habilitadas para Lake Formation. Como o driver do sistema é executado com permissões elevadas, os eventos e logs que o driver do sistema gera podem incluir informações confidenciais. Para evitar que usuários ou códigos não autorizados acessem esses dados confidenciais, o EMR Sem Servidor desabilita o acesso aos logs do driver do sistema.
- Os logs do perfil do sistema são sempre mantidos no armazenamento gerenciado: essa é uma configuração obrigatória que não pode ser desabilitada. Esses registros são armazenados com segurança e criptografados usando uma chave KMS gerenciada pelo cliente ou uma chave KMS AWS gerenciada.

Iceberg

Analise as seguintes considerações ao usar o Apache Iceberg:

- Você só pode usar o Apache Iceberg com o catálogo de sessões e não com catálogos nomeados arbitrariamente.
- As tabelas do Iceberg registradas no Lake Formation oferecem suporte apenas às tabelas de metadados `history`, `metadata_log_entries`, `snapshots`, `files`, `manifests` e `refs`. O Amazon EMR oculta as colunas que podem conter dados confidenciais, como `partitions`, `path` e `summaries`. Essa limitação não se aplica às tabelas do Iceberg que não estão registradas no Lake Formation.
- As tabelas não registradas no Lake Formation oferecem suporte a todos os procedimentos armazenados do Iceberg. Os procedimentos `register_table` e `migrate` não são compatíveis com nenhuma tabela.
- Sugerimos que você use o Iceberg `DataFrameWriter V2` em vez do `V1`.

Solução de problemas

Consulte as seções a seguir para soluções de problemas.

Registro em log

O EMR Sem Servidor usa perfis de recursos do Spark para dividir a execução do trabalho. O EMR Sem Servidor usa o perfil do usuário para executar o código fornecido, enquanto o perfil do sistema

impõe as políticas do Lake Formation. Você pode acessar os logs das tarefas executadas como perfil de usuário.

Para obter mais informações sobre depuração de trabalhos habilitados para o Lake Formation, consulte [Depuração de trabalhos](#).

Interface do usuário do Live e servidor de histórico do Spark

A interface de usuário do Live e o servidor de histórico do Spark têm todos os eventos do Spark gerados no perfil do usuário e os eventos editados gerados pelo driver do sistema.

Você pode ver todas as tarefas dos drivers do usuário e do sistema na guia Executores. No entanto, os links de logs estão disponíveis somente para o perfil do usuário. Além disso, algumas informações são retiradas da interface de usuário do Live, como o número de registros de saída.

O trabalho falhou com permissões insuficientes do Lake Formation

Certifique-se de que o perfil de runtime do trabalho tenha as permissões para executar SELECT e DESCRIBE na tabela que você está acessando.

Falha na execução do trabalho com RDD

No momento, o EMR Sem Servidor não oferece suporte a operações de conjunto de dados resilientes distribuídos (RDD) em trabalhos habilitados para o Lake Formation.

Não é possível acessar arquivos de dados no Amazon S3

Certifique-se de ter registrado a localização do data lake no Lake Formation.

Exceção de validação de segurança

O EMR Sem Servidor detectou um erro de validação de segurança. Entre em contato com o AWS suporte para obter assistência.

Compartilhando o catálogo de dados e tabelas do AWS Glue entre contas

Você pode compartilhar bancos de dados e tabelas entre contas e ainda usar o Lake Formation. Para obter mais informações, consulte [Compartilhamento de dados entre contas no Lake Formation](#) e [Como faço para compartilhar o catálogo de dados e tabelas do AWS Glue usando várias contas?](#) AWS Lake Formation.

API de controle de acesso refinado nativo do Spark permitida PySpark

Para manter a segurança e os controles de acesso aos dados, o controle de acesso refinado (FGAC) do Spark restringe determinadas funções. PySpark Essas restrições são aplicadas por meio de:

- Bloqueio explícito que impede a execução da função
- Incompatibilidades de arquitetura que tornam as funções não funcionais
- Funções que podem gerar erros, retornar mensagens de acesso negado ou não fazer nada quando chamadas

Os seguintes PySpark recursos não são compatíveis com o Spark FGAC:

- Operações RDD (bloqueadas com a exceção do Spark) RDDUnsupported
- Spark Connect (não suportado)
- Spark Streaming (não suportado)

Embora tenhamos testado as funções listadas em um ambiente Native Spark FGAC e confirmado que elas funcionam conforme o esperado, nossos testes normalmente abrangem apenas o uso básico de cada API. Funções com vários tipos de entrada ou caminhos lógicos complexos podem ter cenários não testados.

Para qualquer função não listada aqui e que não faça parte claramente das categorias não suportadas acima, recomendamos:

- Testando-os primeiro em um ambiente gama ou em uma implantação em pequena escala
- Verificando seu comportamento antes de usá-los na produção

Note

Se você ver um método de classe listado, mas não sua classe base, o método ainda deve funcionar — significa apenas que não verificamos explicitamente o construtor da classe base.

A PySpark API é organizada em módulos. O suporte geral para métodos em cada módulo está detalhado na tabela abaixo.

Nome do módulo	Status	Observações
pyspark_core	Compatível	Este módulo contém as principais classes do RDD e, na maioria das vezes, essas funções não são suportadas.
pyspark_sql	Compatível	
teste_de_pyspark	Compatível	
pyspark_resource	Compatível	
pyspark streaming	Bloqueado	O uso do streaming está bloqueado no Spark FGAC.
pyspark_mllib	Experimental	Este módulo contém operações de ML baseadas em RDD, e essas funções geralmente não são suportadas. Este módulo não foi totalmente testado.
pyspark_ml	Experimental	Esse módulo contém operações de ML DataFrame baseadas, e essas funções são suportadas principalmente. Este módulo não foi totalmente testado.
pyspark_pandas	Compatível	
pyspark_pandas_slow	Compatível	

Nome do módulo	Status	Observações
pyspark_connect	Bloqueado	O uso do Spark Connect está bloqueado no Spark FGAC.
pyspark_pandas_connect	Bloqueado	O uso do Spark Connect está bloqueado no Spark FGAC.
pyspark_pandas_slow_connect	Bloqueado	O uso do Spark Connect está bloqueado no Spark FGAC.
pyspark_errors	Experimental	Este módulo não foi totalmente testado. Classes de erro personalizadas não podem ser utilizadas.

Lista de permissões da API

Para uma lista que pode ser baixada e mais fácil de pesquisar, um arquivo com os módulos e classes está disponível nas [funções do Python permitidas no](#) FGAC nativo.

Criptografia entre trabalhadores

Com as versões 6.15.0 e superiores do Amazon EMR, habilite a comunicação criptografada por TLS mútua entre trabalhadores nas execuções de trabalhos do Spark. Quando habilitado, o EMR Sem Servidor gera e distribui automaticamente um certificado exclusivo para cada trabalhador provisionado nas execuções de trabalho. Quando esses trabalhadores se comunicam para trocar mensagens de controle ou transferir dados embaralhados, eles estabelecem uma conexão TLS mútua e usam os certificados configurados para verificar a identidade um do outro. Se um trabalhador não conseguir verificar outro certificado, o handshake do TLS falhará e o EMR Sem Servidor interromperá a conexão entre eles.

Se você estiver usando o Lake Formation com o EMR Sem Servidor, a criptografia TLS mútua está habilitada por padrão.

Como habilitar a criptografia TLS mútua no EMR Sem Servidor

Para habilitar a criptografia TLS mútua na aplicação do Spark, defina `spark.ssl.internode.enabled` como verdadeiro ao [criar a aplicação do EMR Sem Servidor](#). Se você estiver usando o AWS console para criar um aplicativo EMR Serverless, escolha Usar configurações personalizadas, expanda Configuração do aplicativo e insira seu `runtimeConfiguration`

```
aws emr-serverless create-application \  
--release-label emr-6.15.0 \  
--runtime-configuration '{  
  "classification": "spark-defaults",  
  "properties": {"spark.ssl.internode.enabled": "true"}  
}' \  
--type "SPARK"
```

Se você quiser habilitar a criptografia TLS mútua para execuções individuais de trabalhos do Spark, defina `spark.ssl.internode.enabled` como verdadeiro ao usar `spark-submit`.

```
--conf spark.ssl.internode.enabled=true
```

Criptografia de disco com KMS CMK

O EMR Serverless criptografa todos os discos conectados aos trabalhadores por padrão usando chaves de criptografia de propriedade do serviço. Opcionalmente, você pode optar por criptografar esses discos usando suas próprias chaves gerenciadas pelo AWS KMS cliente (). CMKs Isso fornece mais controle sobre suas chaves de criptografia, incluindo a capacidade de estabelecer e manter políticas de chaves e auditar o uso de chaves.

Você pode configurar a criptografia de disco ao criar um aplicativo ou ao enviar trabalhos individuais. Quando ativadas no nível do aplicativo, todas as tarefas desse aplicativo herdam as configurações de criptografia. Você também pode substituir o padrão do aplicativo especificando uma configuração de criptografia de disco ao enviar um trabalho.

Note

A criptografia de disco do EMR Serverless suporta somente chaves KMS simétricas. As chaves do KMS assimétricas não são aceitas. Você deve usar uma chave KMS de

criptografia simétrica que foi criada em. AWS KMS Para obter mais informações sobre AWS KMS, consulte [O que é AWS KMS?](#)

Usando o contexto de criptografia

Opcionalmente, o EMR Serverless usa contexto de criptografia para fornecer dados autenticados adicionais para operações de criptografia. O contexto de criptografia é um conjunto de pares de valores-chave que podem conter dados autenticados adicionais não secretos. O contexto de criptografia é vinculado criptograficamente aos dados criptografados, portanto, o mesmo contexto de criptografia é necessário para descriptografar os dados.

No EMR Serverless, você pode especificar o contexto de criptografia personalizado ao configurar a criptografia de disco. Esse contexto de criptografia é incluído nos AWS CloudTrail registros para ajudar você a identificar e entender suas operações do KMS.

Note

Não armazene informações confidenciais no contexto de criptografia, pois elas aparecem em texto simples nos AWS CloudTrail registros.

Configurando a criptografia de disco com chaves gerenciadas pelo cliente

CreateApplication

Para criptografar discos com sua própria chave KMS, inclua o `diskEncryptionConfiguration` parâmetro ao criar um aplicativo EMR Serverless.

```
aws emr-serverless create-application \  
  --type TYPE \  
  --name APPLICATION_ID \  
  --release-label RELEASE_LABEL \  
  --region AWS_REGION \  
  --disk-encryption-configuration '{  
    "encryptionKeyArn": "key-arn",  
    "encryptionContext": {  
      "key": "value"  
    }  
  }'  
'
```

UpdateApplication

Para atualizar o contexto de criptografia and/or ARN da chave KMS, especifique `diskEncryptionConfiguration` o parâmetro com os novos valores ao atualizar um aplicativo.

```
aws emr-serverless update-application \
  --name APPLICATION_ID \
  --region AWS_REGION \
  --disk-encryption-configuration '{
    "encryptionKeyArn": "key-arn",
    "encryptionContext": {
      "key": "value"
    }
  }'
```

Note

Para cancelar a configuração da criptografia de disco em um aplicativo, passe um vazio `diskEncryptionConfiguration` durante a atualização do aplicativo.

StartJobRun

Para criptografar discos com sua própria chave KMS, use a `diskEncryptionConfiguration` configuração ao enviar uma execução de trabalho.

```
--configuration-overrides '{
  "diskEncryptionConfiguration": {
    "encryptionKeyArn": "key-arn",
    "encryptionContext": {
      "key": "value"
    }
  }
}'
```

Endpoints públicos do Livy

Para criptografar discos com sua própria chave KMS ao criar sessões do Spark por meio de endpoints públicos do Livy, especifique a configuração de criptografia no objeto da sessão. `conf`

```
data = {
```

```

"kind": "pyspark",
"heartbeatTimeoutInSeconds": 60,
"conf": {
  "emr-serverless.session.executionRoleArn": "role_arn",
  "spark.emr-serverless.disk.encryptedKeyArn": "key-arn",
  "spark.emr-serverless.disk.encryptedContext": "key1:value1,key2:value2" #
Optional
}
}

# Send request to create a session with the Livy API endpoint
request = AWSRequest(method='POST', url=endpoint + "/sessions", data=json.dumps(data),
headers=headers)

```

Permissões necessárias para criptografia de disco

Permissões de chave de criptografia para EMR Serverless

Ao criptografar discos com sua própria chave de criptografia, você deve configurar as seguintes permissões de chave KMS para o principal: `emr-serverless.amazonaws.com`

- `kms:GenerateDataKey`: Para gerar chaves de dados para criptografar volumes de disco
- `kms:Decrypt`: Para descriptografar chaves de dados ao acessar o conteúdo do disco criptografado

```

{
  "Effect": "Allow",
  "Principal": {
    "Service": "emr-serverless.amazonaws.com"
  },
  "Action": [
    "kms:Decrypt",
    "kms:GenerateDataKey"
  ],
  "Resource": "*",
  "Condition": {
    "StringLike": {
      "aws:SourceArn": "arn:aws:emr-serverless:region:aws-account-id:/
applications/application-id"
    },
    "StringEquals": {

```

```

        "kms:EncryptionContext:applicationId": "application-id",
        "aws:SourceAccount": "aws-account-id"
    }
}
}

```

Como uma prática recomendada de segurança, adicione uma chave de condição `aws:SourceArn` à política de chave do KMS. A chave de condição global do IAM `aws:SourceArn` ajuda a garantir que o EMR Sem Servidor use a chave do KMS somente para o ARN da aplicação. Além disso, incluir a chave de `aws:SourceAccount` condição fornece outra camada de segurança ao restringir o uso da sua chave KMS às solicitações originadas do ID da AWS conta especificado na condição.

O perfil de runtime do trabalho deve ter as seguintes permissões na política do IAM:

```

{
  "Sid": "Enable GDK and Decrypt",
  "Version": "2012-10-17",
  "Statement": {
    "Effect": "Allow",
    "Action": [
      "kms:GenerateDataKey",
      "kms:Decrypt"
    ],
    "Resource": "key-arn"
  }
}

```

Permissões obrigatórias do usuário

O usuário que envia o trabalho deve ter permissões para usar a chave. Você pode especificar as permissões na política de chave do KMS ou na política do IAM para o usuário, grupo ou perfil. Se o usuário que envia o trabalho não tiver as permissões da chave do KMS, o EMR Sem Servidor rejeitará o envio da execução do trabalho.

Exemplo de política de chaves

A política de chaves a seguir fornece as permissões para `kms:DescribeKey`, `kms:GenerateDataKey` e `kms:Decrypt`:

- `kms:DescribeKey`: para verificar se a chave KMS gerenciada pelo cliente está ativada e é SIMÉTRICA antes de usá-la.

```

{
  "Sid": "Enable DescribeKey",
  "Effect": "Allow",
  "Principal":{
    "AWS": "arn:aws:iam::111122223333:user/user-name"
  },
  "Action": [
    "kms:DescribeKey"
  ],
  "Resource": "*"
},
{
  "Sid": "Enable GDK and Decrypt",
  "Effect": "Allow",
  "Principal":{
    "AWS": "arn:aws:iam::111122223333:user/user-name"
  },
  "Action": [
    "kms:GenerateDataKey",
    "kms:Decrypt"
  ],
  "Resource": "*",
  "Condition": {
    "StringEquals": {
      "kms:ViaService": "emr-serverless.region.amazonaws.com",
      "kms:EncryptionContext:key": "value"
    }
  }
}
}

```

Como uma prática recomendada de segurança, adicione uma chave de condição `kms:viaService` à política de chave do KMS. Isso limita o uso da chave KMS para solicitações de validação apenas do `emr-serverless`.

Exemplo de política do IAM

A política do IAM a seguir fornece as permissões para `kms:DescribeKey` `kms:GenerateDataKey` `kms:Decrypt` e.

```

{
  "Version": "2012-10-17",
  "Statement": {
    "Effect": "Allow",

```

```
    "Action": [
      "kms:DescribeKey",
      "kms:GenerateDataKey",
      "kms:Decrypt"
    ],
    "Resource": "key-arn"
  }
}
```

Monitorando o uso da chave

Você pode monitorar o uso de suas chaves gerenciadas pelo cliente no EMR Serverless por meio de AWS CloudTrail. AWS CloudTrail captura todas as chamadas de API para AWS KMS como eventos, incluindo chamadas do console do EMR Serverless, da API do EMR Serverless, da CLI ou do SDK.

AWS AWS

As informações capturadas incluem o contexto de criptografia que você especificou, o que pode ajudá-lo a identificar e auditar os recursos específicos do EMR Serverless que usaram sua chave KMS. Por exemplo, você pode ver eventos semelhantes aos seguintes em AWS CloudTrail. Para obter mais informações sobre o uso AWS CloudTrail, consulte o [Guia AWS CloudTrail do usuário](#).

GenerateDataKey

Exemplo de evento para GenerateDataKey operações quando o EMR Serverless está criando volumes de disco criptografados

```
{
  "eventVersion": "1.11",
  "userIdentity": {
    "type": "AWSService",
    "principalId": "user",
    "invokedBy": "AWS Internal"
  },
  "eventTime": "2025-07-28T21:43:51Z",
  "eventSource": "kms.amazonaws.com",
  "eventName": "GenerateDataKey",
  "awsRegion": "us-west-2",
  "sourceIPAddress": "ipAddress",
  "userAgent": "userAgent",
  "requestParameters": {
    "encryptionContext": {
      "applicationId": "test"
    }
  }
}
```

```

    },
    "keyId": "arn:aws:kms:region:accountId:key/ffffffff-fffff-aaaaa-eeee-sample",
    "keySpec": "AES_256"
  },
  "responseElements": null,
  "additionalEventData": {
    "keyMaterialId":
"145c963debe558dfb01848d2a4539da940f3478852f86cfe2f52d5df796a5a02"
  },
  "requestID": "cc9d1c5e-97c4-4a4f-ae7a-e576sample",
  "eventID": "0b0fef09-f28d-4da8-a5a1-17b74sample",
  "readOnly": true,
  "resources": [
    {
      "accountId": "account",
      "type": "AWS::KMS::Key",
      "ARN": "arn:aws:kms:region:accountId:key/ffffffff-fffff-aaaaa-eeee-sample"
    }
  ],
  "eventType": "AwsApiCall",
  "managementEvent": true,
  "recipientAccountId": "accountId",
  "eventCategory": "Management"
}

```

Decrypt

Exemplo de evento para operações de descryptografia quando o EMR Serverless está acessando dados criptografados.

```

{
  "eventVersion": "1.11",
  "userIdentity": {
    "type": "AWSService",
    "principalId": "user",
    "invokedBy": "AWS Internal"
  },
  "eventTime": "2025-07-28T21:43:51Z",
  "eventSource": "kms.amazonaws.com",
  "eventName": "Decrypt",
  "awsRegion": "us-west-2",
  "sourceIPAddress": "ipAddress",
  "userAgent": "userAgent",

```

```
"requestParameters": {
  "encryptionContext": {
    "applicationId": "test"
  },
  "keyId": "arn:aws:kms:region:accountId:key/ffffffff-fffff-aaaaa-eeee-sample",
  "keySpec": "AES_256"
},
"responseElements": null,
"additionalEventData": {
  "keyMaterialId":
"145c963debe558dfb01848d2a4539da940f3478852f86cfe2f52d5df796a5a02"
},
"requestID": "cc9d1c5e-97c4-4a4f-ae7a-e576sample",
"eventID": "0b0fef09-f28d-4da8-a5a1-17b74sample",
"readOnly": true,
"resources": [
  {
    "accountId": "account",
    "type": "AWS::KMS::Key",
    "ARN": "arn:aws:kms:region:accountId:key/ffffffff-fffff-aaaaa-eeee-sample"
  }
],
"eventType": "AwsApiCall",
"managementEvent": true,
"recipientAccountId": "accountId",
"eventCategory": "Management"
}
```

Saiba mais

Os recursos a seguir fornecem mais informações sobre a criptografia de dados em repouso.

- Para obter mais informações sobre conceitos AWS KMS básicos, consulte o [Guia do AWS KMS desenvolvedor](#).
- Para obter mais informações sobre as melhores práticas de segurança para AWS KMS, consulte o [Guia do AWS KMS desenvolvedor](#).

Secrets Manager para proteção de dados com o EMR Sem Servidor

AWS Secrets Manager é um serviço de armazenamento secreto para proteger credenciais de banco de dados, chaves de API e outras informações secretas. Em seguida, no seu código, substitua credenciais codificadas por uma chamada de API para o Secrets Manager. Isso ajuda a garantir que o segredo não possa ser comprometido por alguém que esteja examinando seu código, pois o segredo não está ali. Para obter uma visão geral, consulte o [Guia do usuário do AWS Secrets Manager](#).

O Secrets Manager criptografa segredos usando AWS Key Management Service chaves. Para obter mais informações, consulte [Criptografia e descriptografia de segredos](#) no Guia do usuário do AWS Secrets Manager .

Você pode configurar o Secrets Manager para alterar automaticamente os segredos para você de acordo com a programação que você especificar. Isso permite substituir segredos de longo prazo por outros de curto prazo, ajudando a reduzir de maneira significativa o risco de comprometimento. Para obter mais informações, consulte [Alternar os segredos do AWS Secrets Manager](#) no Guia do usuário do AWS Secrets Manager .

O Amazon EMR Serverless se integra AWS Secrets Manager para que você possa armazenar seus dados no Secrets Manager e usar o ID secreto em suas configurações.

Como o EMR Sem Servidor usa segredos

Quando você armazena seus dados no Secrets Manager e usa o ID secreto em suas configurações para o EMR Serverless, você não passa dados de configuração confidenciais para o EMR Serverless em texto simples e os expõe ao ambiente externo. APIs Se você indicar que um par de chave-valor contém o ID do segredo armazenado no Secrets Manager, o EMR Sem Servidor recuperará o segredo ao enviar dados de configuração aos trabalhadores para execução de trabalhos.

Para indicar que um par de chave-valor de uma configuração contém uma referência a um segredo armazenado no Secrets Manager, adicione a anotação `EMR.secret@` ao valor da configuração. Para qualquer propriedade de configuração com anotação de ID secreta, o EMR Sem Servidor chama o Secrets Manager e resolve o segredo no momento da execução do trabalho.

Como criar um segredo

Para criar um segredo, siga as etapas em [Criar um AWS Secrets Manager segredo](#) no Guia do AWS Secrets Manager usuário. Na Etapa 3, escolha o campo Texto sem formatação para inserir o valor sigiloso.

Fornecimento de um segredo em uma classificação de configuração

Os exemplos a seguir mostram como fornecer um segredo em uma classificação de configuração em StartJobRun. Se quiser configurar classificações para o Secrets Manager no nível da aplicação, consulte [Configuração padrão de aplicações do EMR Sem Servidor](#).

Nos exemplos, substitua *SecretName* pelo nome do segredo a ser recuperado. Para obter mais informações, consulte [Como criar um segredo](#).

Nesta seção

- [Especifique referências secretas: Spark](#)
- [Especifique referências secretas: Hive](#)

Especifique referências secretas: Spark

Example— Especifique referências secretas na configuração externa da metastore do Hive para o Spark

```
aws emr-serverless start-job-run \  
  --application-id "application-id" \  
  --execution-role-arn "job-role-arn" \  
  --job-driver '{  
    "sparkSubmit": {  
      "entryPoint": "s3://amzn-s3-demo-bucket/scripts/spark-jdbc.py",  
      "sparkSubmitParameters": "--jars s3://amzn-s3-demo-bucket/mariadb-  
connector-java.jar  
      --conf  
spark.hadoop.java.x.jdo.option.ConnectionDriverName=org.mariadb.jdbc.Driver  
      --conf spark.hadoop.java.x.jdo.option.ConnectionUserName=connection-user-  
name  
      --conf  
spark.hadoop.java.x.jdo.option.ConnectionPassword=EMR.secret@SecretName  
      --conf spark.hadoop.java.x.jdo.option.ConnectionURL=jdbc:mysql://db-host:db-  
port/db-name
```

```

        --conf spark.driver.cores=2
        --conf spark.executor.memory=10G
        --conf spark.driver.memory=6G
        --conf spark.executor.cores=4"
    }
}' \
--configuration-overrides '{
    "monitoringConfiguration": {
        "s3MonitoringConfiguration": {
            "logUri": "s3://amzn-s3-demo-bucket/spark/logs/"
        }
    }
}'

```

Example— Especifique referências secretas para a configuração externa da metastore do Hive na classificação spark-defaults.

```

{
    "classification": "spark-defaults",
    "properties": {

        "spark.hadoop.javax.jdo.option.ConnectionDriverName": "org.mariadb.jdbc.Driver"
        "spark.hadoop.javax.jdo.option.ConnectionURL": "jdbc:mysql://db-host:db-  
port/db-name"
        "spark.hadoop.javax.jdo.option.ConnectionUserName": "connection-user-name"
        "spark.hadoop.javax.jdo.option.ConnectionPassword":
        "EMR.secret@SecretName",
    }
}

```

Especifique referências secretas: Hive

Example— Especifique referências secretas na configuração externa da metastore do Hive para o Hive

```

aws emr-serverless start-job-run \
  --application-id "application-id" \
  --execution-role-arn "job-role-arn" \
  --job-driver '{
    "hive": {
        "query": "s3://amzn-s3-demo-bucket/emr-serverless-hive/query/hive-query.q1",

```

```

    "parameters": "--hiveconf hive.exec.scratchdir=s3://amzn-s3-demo-bucket/emr-
serverless-hive/hive/scratch
                --hiveconf hive.metastore.warehouse.dir=s3://amzn-s3-demo-bucket/
emr-serverless-hive/hive/warehouse
                --hiveconf javax.jdo.option.ConnectionUserName=username
                --hiveconf
javax.jdo.option.ConnectionPassword=EMR.secret@SecretName
                --hiveconf
hive.metastore.client.factory.class=org.apache.hadoop.hive.ql.metadata.SessionHiveMetaStoreCli
                --hiveconf
javax.jdo.option.ConnectionDriverName=org.mariadb.jdbc.Driver
                --hiveconf javax.jdo.option.ConnectionURL=jdbc:mysql://db-host:db-
port/db-name"
    }
}' \
--configuration-overrides '{
  "monitoringConfiguration": {
    "s3MonitoringConfiguration": {
      "logUri": "s3://amzn-s3-demo-bucket"
    }
  }
}'
}'

```

Example— Especifique referências secretas para a configuração externa da metastore do Hive na classificação hive-site.

```

{
  "classification": "hive-site",
  "properties": {
    "hive.metastore.client.factory.class":
"org.apache.hadoop.hive.ql.metadata.SessionHiveMetaStoreClientFactory",
    "javax.jdo.option.ConnectionDriverName": "org.mariadb.jdbc.Driver",
    "javax.jdo.option.ConnectionURL": "jdbc:mysql://db-host:db-port/db-name",
    "javax.jdo.option.ConnectionUserName": "username",
    "javax.jdo.option.ConnectionPassword": "EMR.secret@SecretName"
  }
}

```

Concessão de acesso ao EMR Sem Servidor para recuperar o segredo

Para permitir que o EMR Sem Servidor recupere o valor do segredo do Secrets Manager, adicione a instrução de política a seguir ao seu segredo ao criá-lo. Você deve criar o segredo com a chave

do KMS gerenciada pelo cliente para que o EMR Sem Servidor leia o valor do segredo. Para obter mais informações, consulte [Permissões para a chave do KMS](#) no Guia do usuário do AWS Secrets Manager .

Na política a seguir, substitua *applicationId* pelo ID da aplicação.

Política de recursos do segredo

Você deve incluir as permissões a seguir na política de recursos do segredo no AWS Secrets Manager para permitir que o EMR Sem Servidor recupere valores do segredo. Para garantir que somente uma aplicação específica possa recuperar esse segredo, você pode, opcionalmente, especificar o ID da aplicação do EMR Sem Servidor como uma condição na política.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "secretsmanager:GetSecretValue",
        "secretsmanager:DescribeSecret"
      ],
      "Resource": [
        "*"
      ],
      "Condition": {
        "StringEquals": {
          "aws:SourceArn": "arn:aws:emr-serverless:*:123456789012:/applications/"
        }
      }
    }
  ],
  "Sid": "AllowSECRETSMANAGERGetsecretvalue"
}
```

Crie seu segredo com a seguinte política para a chave gerenciada pelo cliente AWS Key Management Service (AWS KMS):

Política para chave gerenciada pelo cliente AWS KMS

```
{
  "Sid": "Allow EMR Serverless to use the key for decrypting secrets",
  "Effect": "Allow",
  "Principal": {
    "Service": [
      "emr-serverless.amazonaws.com"
    ]
  },
  "Action": [
    "kms:Decrypt",
    "kms:DescribeKey"
  ],
  "Resource": "*",
  "Condition": {
    "StringEquals": {
      "kms:ViaService": "secretsmanager.Região da AWS.amazonaws.com"
    }
  }
}
```

Alternância do segredo

A alternância é quando você atualiza periodicamente um segredo. É possível configurar o AWS Secrets Manager para alternar automaticamente o segredo de acordo com uma programação que você especificar. Dessa forma, é possível substituir segredos de longo prazo por segredos de curto prazo. Isso ajuda a reduzir o risco de comprometimento. O EMR Sem Servidor recupera o valor do segredo de uma configuração anotada quando o trabalho passa para um estado de execução. Se você ou um processo atualizar o valor do segredo no Secrets Manager, você deverá enviar um novo trabalho para que ele possa buscar o valor atualizado.

Note

Os trabalhos que já estão em execução não podem buscar um valor de segredo atualizado. Isso pode resultar em falha no trabalho.

Uso da Concessão de Acesso do Amazon S3 com o EMR Sem Servidor

Visão geral da Concessão de Acesso do S3 para o EMR Sem Servidor

Com as versões 6.15.0 e superiores do Amazon EMR, a Concessão de Acesso do Amazon S3 oferece uma solução de controle de acesso escalável para aumentar o acesso aos seus dados do Amazon S3 a partir do EMR Sem Servidor. Se você tiver uma configuração de permissão complexa ou grande para os dados do S3, use a funcionalidade Concessão de Acesso para escalar as permissões de dados do S3 para usuários, perfis e aplicações.

Use a Concessão de Acesso do S3 para aumentar o acesso aos dados do Amazon S3 além das permissões concedidas pelo perfil de runtime ou pelos perfis do IAM que estão anexados às identidades com acesso à aplicação do EMR Sem Servidor.

Para obter mais informações, consulte [Gerenciar o acesso com a Concessão de Acesso do S3 para Amazon EMR](#), no Guia de gerenciamento do Amazon EMR, e [Gerenciar o acesso com a Concessão de Acesso do S3](#), no Guia do usuário do Amazon Simple Storage Service.

Esta seção descreve como iniciar uma aplicação do EMR Sem Servidor que usa a Concessão de Acesso do S3 para fornecer acesso aos dados no Amazon S3. Para ver as etapas de uso do S3 Access Grants com outras implantações do Amazon EMR, consulte a seguinte documentação:

- [Using S3 Access Grants with Amazon EMR](#)
- [Using S3 Access Grants with Amazon EMR on EKS](#)

Launch an EMR Serverless application with S3 Access Grants for data management

Você pode habilitar a Concessão de Acesso do S3 no EMR Sem Servidor e executar uma aplicação do Spark. Quando sua aplicação solicita dados do S3, o Amazon S3 fornece credenciais temporárias que têm como escopo o bucket, prefixo ou objeto específico.

1. Configure um perfil de execução de trabalhos para a aplicação do EMR Sem Servidor. Inclua as permissões do IAM necessárias para executar os trabalhos do Spark, e use a Concessão de Acesso do S3, `s3:GetDataAccess` e `s3:GetAccessGrantsInstanceForPrefix`:

```
{
```

```

"Effect": "Allow",
"Action": [
"s3:GetDataAccess",
"s3:GetAccessGrantsInstanceForPrefix"
],
"Resource": [ //LIST ALL INSTANCE ARNS THAT THE ROLE IS ALLOWED TO QUERY
  "arn:aws_partition:s3:Region:account-id1:access-grants/default",
  "arn:aws_partition:s3:Region:account-id2:access-grants/default"
]
}

```

Note

Se você especificar perfis do IAM para a execução de trabalhos que tenham permissões adicionais de acesso direto ao S3, os usuários poderão acessar os dados permitidos pelo perfil mesmo que não tenham permissão da Concessão de Acesso do S3.

2. Inicie a aplicação do EMR Sem Servidor com um rótulo de lançamento 6.15.0 ou superior do Amazon EMR e a classificação `spark-defaults`, conforme mostra o exemplo a seguir. Substitua os valores em *red text* pelos valores apropriados ao seu cenário de uso.

```

aws emr-serverless start-job-run \
  --application-id application-id \
  --execution-role-arn job-role-arn \
  --job-driver '{
    "sparkSubmit": {
      "entryPoint": "s3://us-east-1.elasticmapreduce/emr-containers/samples/
wordcount/scripts/wordcount.py",
      "entryPointArguments": ["s3://amzn-s3-demo-destination-bucket1/
wordcount_output"],
      "sparkSubmitParameters": "--conf spark.executor.cores=1 --conf
spark.executor.memory=4g --conf spark.driver.cores=1 --conf spark.driver.memory=4g
--conf spark.executor.instances=1"
    }
  }' \
  --configuration-overrides '{
    "applicationConfiguration": [{
      "classification": "spark-defaults",
      "properties": {
        "spark.hadoop.fs.s3.s3AccessGrants.enabled": "true",
        "spark.hadoop.fs.s3.s3AccessGrants.fallbackToIAM": "false"
      }
    }
  ]
}

```

```
} ]  
' }
```

Considerações sobre a Concessão de Acesso do S3 com o EMR Sem Servidor

Para obter informações importantes sobre suporte, compatibilidade e comportamento ao usar a Concessão de Acesso do Amazon S3 com o EMR Sem Servidor, consulte [Considerações sobre a Concessão de Acesso do S3 com o Amazon EMR](#) no Guia de gerenciamento do Amazon EMR.

Registro de chamadas de API sem servidor do Amazon EMR usando AWS CloudTrail

O Amazon EMR Serverless é integrado com AWS CloudTrail, um serviço que fornece um registro das ações realizadas por um usuário, função ou serviço AWS no EMR Serverless. CloudTrail captura todas as chamadas de API para o EMR Serverless como eventos. As chamadas capturadas incluem chamadas ao console do EMR Sem Servidor e as chamadas ao código para as operações da API do EMR Sem Servidor. Se você criar uma trilha, habilite a entrega contínua de CloudTrail eventos para um bucket do Amazon S3, incluindo eventos para o EMR Serverless. Se você não configurar uma trilha, ainda poderá acessar os eventos mais recentes no CloudTrail console no Histórico de eventos. Usando as informações coletadas por CloudTrail, você pode determinar a solicitação que foi feita ao EMR Serverless, o endereço IP do qual a solicitação foi feita, quem fez a solicitação, quando ela foi feita e detalhes adicionais.

Para saber mais, consulte o [Guia do usuário do AWS CloudTrail](#).

Informações do EMR Serverless em CloudTrail

CloudTrail é ativado no seu Conta da AWS quando você cria a conta. Quando a atividade ocorre no EMR Serverless, essa atividade é registrada em um CloudTrail evento junto com outros eventos de AWS serviço no histórico de eventos. Você pode acessar, pesquisar e baixar eventos recentes na sua Conta da AWS. Para obter mais informações, consulte [Visualização de eventos com histórico de CloudTrail eventos](#).

Para um registro contínuo dos eventos em seu Conta da AWS, incluindo eventos para o EMR Serverless, crie uma trilha. Uma trilha permite CloudTrail entregar arquivos de log para um bucket do Amazon S3. Por padrão, quando você cria uma trilha no console, ela é aplicada a todas as Regiões

da AWS. A trilha registra eventos de todas as regiões na AWS partição e entrega os arquivos de log ao bucket do Amazon S3 que você especificar. Além disso, configure outros AWS serviços para analisar e agir com base nos dados de eventos coletados nos CloudTrail registros. Para obter mais informações, consulte o seguinte:

- [Visão geral da criação de uma trilha](#)
- [CloudTrail serviços e integrações suportados](#)
- [Configurando notificações do Amazon SNS para CloudTrail](#)
- [Recebendo arquivos de CloudTrail log de várias regiões](#) e [Recebendo arquivos de CloudTrail log de várias contas](#)

[Todas as ações do EMR Serverless são registradas CloudTrail e documentadas na Referência da API do EMR Serverless.](#) Por exemplo, chamadas para o `CreateApplication StartJobRun` e `CancelJobRun` as ações geram entradas nos arquivos de CloudTrail log.

Cada entrada de log ou evento contém informações sobre quem gerou a solicitação. As informações de identidade ajudam a determinar o seguinte:

- Se a solicitação foi feita com credenciais de usuário root ou AWS Identity and Access Management (IAM).
- Se a solicitação foi feita com credenciais de segurança temporárias de um perfil ou de um usuário federado.
- Se a solicitação foi feita por outro AWS serviço.

Para obter mais informações, consulte o elemento [CloudTrail userIdentity](#).

Noções básicas das entradas do arquivo de log do EMR Sem Servidor

Uma trilha é uma configuração que permite a entrega de eventos como arquivos de log para um bucket do Amazon S3 que você especificar. CloudTrail os arquivos de log contêm uma ou mais entradas de log. Um evento representa uma única solicitação de qualquer fonte e inclui informações sobre a ação solicitada, a data e a hora da ação, os parâmetros da solicitação e assim por diante. CloudTrail os arquivos de log não são um rastreamento de pilha ordenado das chamadas públicas de API, portanto, eles não aparecem em nenhuma ordem específica.

O exemplo a seguir mostra uma entrada de CloudTrail registro que demonstra a `CreateApplication` ação.

```

{
  "eventVersion": "1.08",
  "userIdentity": {
    "type": "AssumedRole",
    "principalId": "AIDACKCEVSQ6C2EXAMPLE:admin",
    "arn": "arn:aws:sts::012345678910:assumed-role/Admin/admin",
    "accountId": "012345678910",
    "accessKeyId": "AKIAIOSFODNN7EXAMPLE",
    "sessionContext": {
      "sessionIssuer": {
        "type": "Role",
        "principalId": "AIDACKCEVSQ6C2EXAMPLE",
        "arn": "arn:aws:iam::012345678910:role/Admin",
        "accountId": "012345678910",
        "userName": "Admin"
      },
      "webIdFederationData": {},
      "attributes": {
        "creationDate": "2022-06-01T23:46:52Z",
        "mfaAuthenticated": "false"
      }
    }
  },
  "eventTime": "2022-06-01T23:49:28Z",
  "eventSource": "emr-serverless.amazonaws.com",
  "eventName": "CreateApplication",
  "awsRegion": "us-west-2",
  "sourceIPAddress": "203.0.113.0",
  "userAgent": "PostmanRuntime/7.26.10",
  "requestParameters": {
    "name": "my-serverless-application",
    "releaseLabel": "emr-6.6",
    "type": "SPARK",
    "clientToken": "0a1b234c-de56-7890-1234-567890123456"
  },
  "responseElements": {
    "name": "my-serverless-application",
    "applicationId": "1234567890abcdef0",
    "arn": "arn:aws:emr-serverless:us-west-2:555555555555:/
applications/1234567890abcdef0"
  },
  "requestID": "890b8639-e51f-11e7-b038-EXAMPLE",
  "eventID": "874f89fa-70fc-4798-bc00-EXAMPLE",

```

```
"readOnly": false,  
"eventType": "AwsApiCall",  
"managementEvent": true,  
"recipientAccountId": "012345678910",  
"eventCategory": "Management"  
}
```

Validação de conformidade do Amazon EMR Sem Servidor

A segurança e a conformidade do EMR Serverless são avaliadas por auditores terceirizados como parte de vários programas de AWS conformidade, incluindo os seguintes:

- Controles do Sistema e da Organização (CSO)
- Padrão de segurança de dados do setor de cartão de pagamento (PCI DSS – Payment Card Industry Data Security Standard)
- Federal Risk and Authorization Management Program (FedRAMP) Moderate
- Health Insurance Portability and Accountability Act (HIPAA)

AWS fornece uma lista frequentemente atualizada de AWS serviços no escopo de programas de conformidade específicos em [AWS Serviços no escopo por programa de conformidade](#).

Relatórios de auditoria de terceiros estão disponíveis para você baixar usando AWS Artifact. Para obter mais informações, consulte [Baixar relatórios no AWS Artifact](#).

Para obter mais informações sobre programas de AWS conformidade, consulte [Programas de AWS conformidade](#).

Sua responsabilidade de conformidade ao usar o EMR Sem Servidor é determinada pela confidencialidade de seus dados, pelas metas de conformidade da sua empresa e pelas regulamentações e leis aplicáveis. Caso seu uso do EMR Sem Servidor esteja sujeito à conformidade com padrões como HIPAA, PCI ou FedRAMP Moderate, a AWS fornecerá recursos para ajudar:

- [Guias de início rápido sobre segurança e conformidade](#) que discutem considerações arquitetônicas e etapas para a implantação de ambientes básicos focados em segurança e conformidade em AWS
- [AWS Os guias de conformidade do cliente](#) podem ajudá-lo a entender o modelo de responsabilidade compartilhada sob a ótica da conformidade. Os guias resumem as práticas

recomendadas para proteção de Serviços da AWS e mapeiam as diretrizes para controles de segurança em várias estruturas (incluindo o Instituto Nacional de Padrões e Tecnologia (NIST), o Conselho de Padrões de Segurança do Setor de Cartões de Pagamento (PCI) e a Organização Internacional de Padronização (ISO)).

- [AWS Config](#) O pode ser usado para avaliar até que ponto suas configurações de recursos atendem adequadamente a práticas internas e a diretrizes e regulamentações do setor.
- [AWS Os Recursos de Conformidade](#) são uma coleção de pastas de trabalho e guias que podem ser aplicados ao seu setor e localização.
- AWS O [Security Hub](#) fornece uma visão abrangente do seu estado de segurança interno AWS e ajuda você a verificar sua conformidade com os padrões e as melhores práticas do setor de segurança.
- [AWS Audit Manager](#)— isso AWS service (Serviço da AWS) ajuda você a auditar continuamente seu AWS uso para simplificar a forma como você gerencia o risco e a conformidade com as regulamentações e os padrões do setor.

Resiliência no Amazon EMR Sem Servidor

A infraestrutura AWS global é construída em torno de AWS regiões e zonas de disponibilidade. AWS As regiões fornecem várias zonas de disponibilidade fisicamente separadas e isoladas, conectadas a redes de baixa latência, alta taxa de transferência e alta redundância. Com zonas de disponibilidade, projete e opere aplicações e bancos de dados que automaticamente executam o failover entre zonas sem interrupção. As zonas de disponibilidade são altamente disponíveis, tolerantes a falhas e escaláveis que uma ou várias infraestruturas de data center tradicionais.

Para obter mais informações sobre AWS regiões e zonas de disponibilidade, consulte [Infraestrutura AWS global](#).

Além da infraestrutura AWS global, o Amazon EMR Serverless oferece integração com o Amazon S3 por meio do EMRFS para ajudar a suportar suas necessidades de resiliência e backup de dados.

Segurança da infraestrutura no Amazon EMR Sem Servidor

Como um serviço gerenciado, o Amazon EMR é protegido pela segurança de rede AWS global. Para obter informações sobre serviços AWS de segurança e como AWS proteger a infraestrutura, consulte [AWS Cloud Security](#). Para projetar seu AWS ambiente usando as melhores práticas de segurança

de infraestrutura, consulte [Proteção](#) de infraestrutura no Security Pillar AWS Well-Architected Framework.

Você usa chamadas de API AWS publicadas para acessar o Amazon EMR pela rede. Os clientes devem oferecer compatibilidade com:

- Transport Layer Security (TLS). Exigimos TLS 1.2 e recomendamos TLS 1.3.
- Conjuntos de criptografia com perfect forward secrecy (PFS) como DHE (Ephemeral Diffie-Hellman) ou ECDHE (Ephemeral Elliptic Curve Diffie-Hellman). A maioria dos sistemas modernos, como Java 7 e versões posteriores, comporta esses modos.

Análise de configuração e vulnerabilidade no Amazon EMR Sem Servidor

AWS lida com tarefas básicas de segurança, como sistema operacional (SO) convidado e aplicação de patches em bancos de dados, configuração de firewall e recuperação de desastres. Esses procedimentos foram revisados e certificados por terceiros certificados. Para obter mais detalhes, consulte os seguintes recursos:

- [Validação de conformidade do Amazon EMR Sem Servidor](#)
- [Modelo de responsabilidade compartilhada](#)
- [Amazon Web Services: visão geral dos processos de segurança](#)

Endpoints e cotas do EMR Serverless

Service endpoints

Para se conectar programaticamente a um AWS service (Serviço da AWS), você usa um endpoint. Um endpoint é o URL do ponto de entrada para um serviço da Web da AWS. Além dos endpoints padrão, alguns Serviços da AWS oferecem AWS endpoints FIPS em regiões selecionadas. A tabela a seguir lista os endpoints de serviço do EMS Sem Servidor. Para obter mais informações, consulte [Endpoints de AWS service \(Serviço da AWS\)](#).

Endpoints do serviço EMR Sem Servidor

Nome da região	Região	Endpoint	Protocolo
Leste dos EUA (Ohio)	us-east-2 (limitado às seguintes zonas de disponibilidade: use2-az1, use2-az2, e use2-az3)	emr-serverless.us-east-2.amazonaws.com	HTTPS
Leste dos EUA (Norte da Virgínia)	us-east-1 (limitado às seguintes zonas de disponibilidade: use1-az1, use1-az2, use1-az4, use1-az5 e use1-az6)	emr-serverless.us-east-1.amazonaws.com emr-serverless-fips.us-east-1.amazonaws.com	HTTPS
Oeste dos EUA (N. da Califórnia)	us-west-1	emr-serverless.us-west-1.amazonaws.com	HTTPS
Oeste dos EUA (Oregon)	us-west-2	emr-serverless.us-	HTTPS

Nome da região	Região	Endpoint	Protocolo
		west-2.amazonaws.com emr-serverless-fips.us-west-2.amazonaws.com	
África (Cidade do Cabo)	af-south-1	emr-serverless.af-south-1.amazonaws.com	HTTPS
Ásia-Pacífico (Hong Kong)	ap-east-1	emr-serverless.ap-east-1.amazonaws.com	HTTPS
Ásia-Pacífico (Jacarta)	ap-southeast-3	emr-serverless.ap-southeast-3.amazonaws.com	HTTPS
Ásia-Pacífico (Melbourne)	ap-southeast-4	emr-serverless.ap-southeast-4.amazonaws.com	HTTPS
Ásia-Pacífico (Malásia)	ap-southeast-5	emr-serverless.ap-southeast-5.amazonaws.com	HTTPS

Nome da região	Região	Endpoint	Protocolo
Ásia-Pacífico (Mumbai)	ap-south-1	emr-serverless.ap-south-1.amazonaws.com	HTTPS
Ásia-Pacífico (Osaka)	ap-northeast-3	emr-serverless.ap-northeast-3.amazonaws.com	HTTPS
Ásia-Pacífico (Seul)	ap-northeast-2	emr-serverless.ap-northeast-2.amazonaws.com	HTTPS
Ásia-Pacífico (Singapura)	ap-southeast-1	emr-serverless.ap-southeast-1.amazonaws.com	HTTPS
Ásia-Pacífico (Sydney)	ap-southeast-2	emr-serverless.ap-southeast-2.amazonaws.com	HTTPS
Ásia-Pacífico (Tóquio)	ap-northeast-1	emr-serverless.ap-northeast-1.amazonaws.com	HTTPS

Nome da região	Região	Endpoint	Protocolo
Canadá (Central)	ca-central-1 (limitado às seguintes zonas de disponibilidade: cac1-az1 e cac1-az2)	emr-serverless.ca-central-1.amazonaws.com	HTTPS
Oeste do Canadá (Calgary)	ca-west-1	emr-serverless.ca-west-1.amazonaws.com	HTTPS
Europa (Frankfurt)	eu-central-1	emr-serverless.eu-central-1.amazonaws.com	HTTPS
Europa (Zurique)	eu-central-2	emr-serverless.eu-central-2.amazonaws.com	HTTPS
Europa (Irlanda)	eu-west-1	emr-serverless.eu-west-1.amazonaws.com	HTTPS
Europa (Londres)	eu-west-2	emr-serverless.eu-west-2.amazonaws.com	HTTPS
Europa (Milão)	eu-south-1	emr-serverless.eu-south-1.amazonaws.com	HTTPS

Nome da região	Região	Endpoint	Protocolo
Europa (Paris)	eu-west-3	emr-serverless.eu-west-3.amazonaws.com	HTTPS
Europa (Espanha)	eu-south-2	emr-serverless.eu-south-2.amazonaws.com	HTTPS
Europa (Estocolmo)	eu-north-1	emr-serverless.eu-north-1.amazonaws.com	HTTPS
Israel (Tel Aviv)	il-central-1	emr-serverless.il-central-1.amazonaws.com	HTTPS
Oriente Médio (Bahrein)	me-south-1	emr-serverless.me-south-1.amazonaws.com	HTTPS
Oriente Médio (Emirados Árabes Unidos)	me-central-1	emr-serverless.me-central-1.amazonaws.com	HTTPS
América do Sul (São Paulo)	sa-east-1	emr-serverless.sa-east-1.amazonaws.com	HTTPS

Nome da região	Região	Endpoint	Protocolo
China (Pequim)	cn-north-1 (limitado às seguintes zonas de disponibilidade: cnn1-az1, cnn1-az2)	emr-serverless.cn-north-1.amazonaws.com.cn	HTTPS
AWS GovCloud (Leste dos EUA)	us-gov-east-1	emr-serverless.us-gov-east-1.amazonaws.com	HTTPS
AWS GovCloud (Oeste dos EUA)	us-gov-west-1	emr-serverless.us-gov-west-1.amazonaws.com	HTTPS

Cotas de serviço

As cotas de serviço, também conhecidas como limites, são o número máximo de recursos ou operações de serviço que você Conta da AWS pode usar. O EMR Sem Servidor coleta métricas de uso da cota de serviço a cada minuto e as publica no namespace AWS/Usage.

Note

AWS As novas contas têm cotas iniciais mais baixas que podem aumentar com o tempo. O Amazon EMR Serverless monitora o uso da conta em cada uma e Região da AWS, em seguida, aumenta automaticamente as cotas com base no seu uso.

A tabela a seguir lista as cotas de serviço do EMR Sem Servidor. Para mais informações, consulte [Cotas de AWS service \(Serviço da AWS\)](#).

Nome	Limite padrão	Ajustável?	Description
Máximo de v simultâneos CPUs por conta	16	Sim	O número máximo de v CPUs que pode ser executado simultaneamente para a conta atual Região da AWS.
Máximo de trabalhos em fila por conta	2000	Sim	O número máximo de trabalhos em fila para a conta na Região da AWS atual.

Limites de API

A seguir, descrevemos os limites de API por região para a Conta da AWS.

Recurso	Cota padrão
ListApplications	10 transações por segundo. Expansão de 50 transações por segundo.
CreateApplication	1 transação por segundo. Expansão de 25 transações por segundo.
DeleteApplication	1 transação por segundo. Expansão de 25 transações por segundo.
GetApplication	10 transações por segundo. Expansão de 50 transações por segundo.
UpdateApplication	1 transação por segundo. Expansão de 25 transações por segundo.
ListJobRuns	1 transação por segundo. Expansão de 25 transações por segundo.

Recurso	Cota padrão
StartJobRun	1 transação por segundo. Expansão de 25 transações por segundo.
GetDashboardForJobRun	1 transação por segundo. Expansão de 2 transações por segundo.
CancelJobRun	1 transação por segundo. Expansão de 25 transações por segundo.
GetJobRun	10 transações por segundo. Expansão de 50 transações por segundo.
StartApplication	1 transação por segundo. Expansão de 25 transações por segundo.
StopApplication	1 transação por segundo. Expansão de 25 transações por segundo.

Outras considerações

A lista a seguir contém outras considerações sobre o EMR Sem Servidor.

- O EMR Serverless está disponível no seguinte: Regiões da AWS
 - Leste dos EUA (Ohio)
 - Leste dos EUA (Norte da Virgínia)
 - Oeste dos EUA (N. da Califórnia)
 - Oeste dos EUA (Oregon)
 - África (Cidade do Cabo)
 - Ásia-Pacífico (Hong Kong)
 - Ásia-Pacífico (Jacarta)
 - Ásia-Pacífico (Mumbai)
 - Ásia-Pacífico (Osaka)
 - Ásia-Pacífico (Seul)
 - Ásia-Pacífico (Singapura)
 - Ásia-Pacífico (Sydney)
 - Ásia-Pacífico (Tóquio)
 - Canadá (Central)
 - Europa (Frankfurt)
 - Europa (Irlanda)
 - Europa (Londres)
 - Europa (Milão)
 - Europe (Paris)
 - Europa (Espanha)
 - Europe (Stockholm)
 - Oriente Médio (Bahrein)
 - Oriente Médio (Emirados Árabes Unidos)
 - América do Sul (São Paulo)
- ~~AWS GovCloud (Leste dos EUA)~~
- AWS GovCloud (Oeste dos EUA)

Para obter uma lista de endpoints associados a essas regiões, consulte [Service endpoints](#).

- O tempo limite padrão para a execução de um trabalho é de 12 horas. Você pode alterar essa configuração com a `executionTimeoutMinutes` propriedade na `startJobRun` API ou no AWS SDK. Você pode definir `executionTimeoutMinutes` como 0 se quiser que a execução do trabalho nunca atinja o tempo limite. Por exemplo, se você tiver uma aplicação de streaming, defina `executionTimeoutMinutes` como 0 para permitir que o trabalho de streaming seja executado continuamente.
- A `billedResourceUtilization` propriedade na `getJobRun` API mostra a vCPU, a memória e o armazenamento agregados AWS que foram cobrados pela execução do trabalho. Os recursos cobrados incluem um uso mínimo de 1 minuto para trabalhadores, além de armazenamento adicional de mais de 20 GB por trabalhador. Esses recursos não incluem o uso de trabalhadores pré-inicializados ociosos.
- Sem conectividade VPC, um trabalho pode acessar alguns AWS service (Serviço da AWS) endpoints no mesmo. Região da AWS Esses serviços incluem Amazon S3, AWS Glue, AWS Lake Formation, Amazon CloudWatch Logs, AWS KMS, AWS Security Token Service Amazon DynamoDB, e. AWS Secrets Manager Você pode habilitar a conectividade VPC para acessar outros Serviços da AWS por meio do [AWS PrivateLink](#), mas não é necessário fazer isso. Para acessar serviços externos, você pode criar a aplicação com uma VPC.
- O EMR Sem Servidor não oferece suporte ao HDFS. Os discos locais dos trabalhadores são um armazenamento temporário que o EMR Sem Servidor usa para embaralhar e processar dados durante a execução do trabalho.

Versões de lançamento do Amazon EMR Sem Servidor

Uma versão do Amazon EMR corresponde a um conjunto de aplicações de código aberto do ecossistema de big data. Cada versão inclui aplicações, componentes e recursos de big data que você seleciona para que o Amazon EMR Sem Servidor implante e configure ao executar o trabalho.

Com as versões 6.6.0 e superiores do Amazon EMR, implante o EMR Sem Servidor. Essa opção de implantação não está disponível em versões anteriores do Amazon EMR. Ao enviar seu trabalho, especifique uma das versões compatíveis a seguir.

Tópicos

- [AWS runtime for Apache Spark\(pré-visualização do emr-spark-8.0\)](#)
- [EMR Serverless7.12.0](#)
- [EMR Serverless7.11.0](#)
- [EMR Serverless7.10.0](#)
- [EMR Serverless7.9.0](#)
- [EMR Serverless7.8.0](#)
- [EMR Serverless7.7.0](#)
- [EMR Serverless7.6.0](#)
- [EMR Serverless7.5.0](#)
- [EMR Serverless7.4.0](#)
- [EMR Serverless7.3.0](#)
- [EMR Serverless7.2.0](#)
- [EMR Serverless7.1.0](#)
- [EMR Serverless7.0.0](#)
- [EMR Serverless6.15.0](#)
- [EMR Serverless6.14.0](#)
- [EMR Serverless6.13.0](#)
- [EMR Serverless6.12.0](#)
- [EMR Serverless6.11.0](#)
- [EMR Serverless6.10.0](#)
- [EMR Serverless6.9.0](#)

- [EMR Serverless6.8.0](#)
- [EMR Serverless6.7.0](#)
- [EMR Serverless6.6.0](#)

AWS runtime for Apache Spark(pré-visualização do emr-spark-8.0)

A tabela a seguir lista as versões do aplicativo disponíveis com AWS runtime for Apache Spark (emr-spark-8.0-preview).

Informações da versão da aplicação

Aplicação	Versão
Spark	4.0.1-amzn-0

AWS runtime for Apache SparkNotas de lançamento (emr-spark-8.0-preview)

- Versão prévia — Esta é uma versão prévia AWS runtime for Apache Spark com o Apache Spark 4.0.1. Essa prévia está disponível somente no EMR Serverless.
- Disponibilidade regional — Essa versão prévia está disponível em todas as AWS regiões onde o EMR Serverless está disponível, exceto nas regiões da China e AWS GovCloud (EUA).
- Informações sobre a versão do aplicativo - Esta versão vem com as seguintes versões do aplicativo:
 - AWS SDK for Java 2.35.5, 1.12.792
 - Python 3.9, 3.11, 3.12
 - Scala 2.13.16
 - AmazonCloudWatchAgent 1.300034.0-amzn-0
 - Delta 4.0.0-amzn-0-spark
 - Iceberg 1.10.0-amzn-spark-0
 - Essa versão vem com o Amazon Corretto 17 (construído no OpenJDK) por padrão para aplicativos compatíveis com o Corretto 17 (JDK 17).
- Limitações da versão prévia - Os seguintes recursos não estão disponíveis nesta versão prévia:

- Recursos interativos e de integração: SageMaker Unified Studio, integração com EMR Studio, Spark Connect, Livy JupyterEnterpriseGateway e não são suportados.
- Formatos de tabela e controle de acesso: Hudi, Delta Universal Format e controle de acesso refinado (FGAC) com filtragem e operadores em nível de linha ou coluna não são suportados. DDL/DML
- Conectores de dados: spark-sql-kinesis, emr-dynamodb e spark-redshift não estão disponíveis.
- Servidor de histórico: o Persistent Spark History Server não está disponível nesta versão prévia. Os usuários ainda podem acessar a interface ativa do Spark para monitorar e depurar trabalhos ativos sem servidor em tempo real.
- Características especializadas: As visualizações materializadas não estão disponíveis.
- Recursos de pré-visualização - Você pode testar os seguintes recursos nesta versão prévia. Essa versão prévia não é recomendada para cargas de trabalho de produção:
 - Características do SQL: modo ANSI SQL com tratamento de tipos mais rigoroso, sintaxe SQL PIPE (|>) para operações de encadeamento, tipo de dados VARIANT para dados JSON semiestruturados, scripts SQL com instruções de fluxo de controle e variáveis de sessão e funções SQL definidas pelo usuário.
 - Aprimoramentos de streaming: API de processamento arbitrário com estado v2 com transformWithState operador, leitor de fonte de dados de estado para estado de streaming consultável (experimental) e armazenamento de estado aprimorado com ponto de verificação aprimorado do changelog do RocksDB.
 - Suporte ao formato de tabela: Apache Iceberg v3 com suporte ao tipo de dados VARIANT, integração com tabelas AWS S3 e Full Table Access (FTA) com AWS Lake Formation tabelas Iceberg, Delta Lake e Hive.
- Documentação adicional - Para obter documentação adicional do Apache Spark, consulte a documentação de lançamento do [Apache Spark 4.0.1](#).

Conceitos básicos

Para começar a usar a versão prévia do Apache Spark 4.0.1, crie um aplicativo EMR Serverless usando a CLI: AWS

```
aws emr-serverless create-application --type spark \  
  --release-label emr-spark-8.0-preview \  
  --region us-east-1 --name spark4-preview
```

EMR Serverless 7.12.0

A tabela a seguir lista as versões do aplicativo disponíveis com a EMR Serverless 7.12.0.

Aplicação	Versão
Apache Spark	3.5.6
Apache Hive	3.1.3
Apache Tez	0.10.2

Notas de versão do EMR Serverless 7.12.0

- Novos recursos
 - Armazenamento sem servidor para EMR Serverless — O Amazon EMR serverless introduz o armazenamento sem servidor, com o EMR versão 7.12 e posterior, que elimina o provisionamento de disco local para cargas de trabalho do Apache Spark. O EMR Serverless gerencia automaticamente a operação intermediária de dados, como o shuffle, sem custos de armazenamento. O armazenamento sem servidor separa o armazenamento da computação, permitindo que o Spark libere os trabalhadores imediatamente quando estão ociosos, em vez de mantê-los ativos para preservar dados temporários. Para saber mais, consulte aws.amazon.com/serverless-storage-for-emr-serverless.
 - Visualizações materializadas do Iceberg - A partir do Amazon EMR 7.12.0, o Amazon EMR Spark oferece suporte à criação e ao gerenciamento de visualizações materializadas do Iceberg (MV)
 - Hudi Full Table Access — A partir do Amazon EMR 7.12.0, o Amazon EMR agora oferece suporte ao controle do Full Table Access (FTA) para o Apache Hudi no Apache Spark com base nas suas políticas definidas em Lake Formation. Esse recurso permite operações de leitura e gravação de suas tarefas do Amazon EMR Spark em tabelas registradas do Lake Formation quando a função de trabalho tem acesso total à tabela.
 - Atualização da versão Iceberg — O Amazon EMR 7.12.0 é compatível com o Apache Iceberg versão 1.10

EMR Serverless7.11.0

A tabela a seguir lista as versões do aplicativo disponíveis com a EMR Serverless 7.11.0.

Aplicação	Versão
Apache Spark	3.5.6
Apache Hive	3.1.3
Apache Tez	0.10.2

Notas de versão do EMR Serverless 7.11.0

- Tempo máximo de execução do trabalho — O valor máximo `executionTimeoutMinutes` em `StartJobRun` ação para trabalhos em LOTE é de 7 dias a partir desta versão. `executionTimeoutMinutes` não pode mais ser definido como, 0 ou seja, sem tempo limite, para execuções de trabalhos em lote.

EMR Serverless7.10.0

A tabela a seguir lista as versões do aplicativo disponíveis com a EMR Serverless 7.10.0.

Aplicação	Versão
Apache Spark	3.5.5
Apache Hive	3.1.3
Apache Tez	0.10.2

Notas de versão do EMR Serverless 7.10.0

- Métricas do EMR Sem Servidor: as métricas de monitoramento são reestruturadas para focar nas dimensões `ApplicationName` e `JobName`. As métricas mais antigas não serão mais atualizadas no futuro. Para obter mais informações, consulte [Monitoramento de aplicações e trabalhos do EMR Sem Servidor](#).

EMR Serverless7.9.0

A tabela a seguir lista as versões do aplicativo disponíveis com a EMR Serverless 7.9.0.

Aplicação	Versão
Apache Spark	3.5.5
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless7.8.0

A tabela a seguir lista as versões do aplicativo disponíveis com a EMR Serverless 7.8.0.

Aplicação	Versão
Apache Spark	3.5.4
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless7.7.0

A tabela a seguir lista as versões do aplicativo disponíveis com a EMR Serverless 7.7.0.

Aplicação	Versão
Apache Spark	3.5.3
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless7.6.0

A tabela a seguir lista as versões do aplicativo disponíveis com a EMR Serverless 7.6.0.

Aplicação	Versão
Apache Spark	3.5.3
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless7.5.0

A tabela a seguir lista as versões do aplicativo disponíveis com a EMR Serverless 7.5.0.

Aplicação	Versão
Apache Spark	3.5.2
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless7.4.0

A tabela a seguir lista as versões do aplicativo disponíveis com a EMR Serverless 7.4.0.

Aplicação	Versão
Apache Spark	3.5.2
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless7.3.0

A tabela a seguir lista as versões das aplicações disponíveis com o EMR Serverless 7.3.0.

Aplicação	Versão
Apache Spark	3.5.1
Apache Hive	3.1.3
Apache Tez	0.10.2

Notas de versão do EMR Serverless 7.3.0

- Simultaneidade de trabalhos e enfileiramento com o EMR Sem Servidor: a simultaneidade e o enfileiramento de trabalhos são habilitados por padrão quando você cria uma nova aplicação do EMR Sem Servidor no Amazon EMR versão 7.3.0 ou superior. Para obter mais informações, consulte [the section called “Simultaneidade e enfileiramento de trabalhos”](#), que detalha como começar com a simultaneidade e filas e também contém uma lista de considerações sobre recursos.

EMR Serverless7.2.0

A tabela a seguir lista as versões das aplicações disponíveis com o EMR Serverless 7.2.0.

Aplicação	Versão
Apache Spark	3.5.1
Apache Hive	3.1.3
Apache Tez	0.10.2

Notas da versão 7.2.0 do EMR Sem Servidor

- Lake Formation com EMR Serverless — agora você pode usar AWS Lake Formation para aplicar controles de acesso refinados em tabelas do catálogo de dados que são apoiadas pelo S3. Esse

recurso permite configurar controles de acesso em nível de tabela, linha, coluna e célula para consultas de leitura nos trabalhos do Spark no EMR Sem Servidor. Para obter mais informações, consulte [the section called “Lake Formation para o FGAC”](#) e [the section called “Considerações e limitações”](#).

EMR Serverless7.1.0

A tabela a seguir lista as versões das aplicações disponíveis com o EMR Serverless 7.1.0.

Aplicação	Versão
Apache Spark	3.5.0
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless7.0.0

A tabela a seguir lista as versões das aplicações disponíveis com o EMR Serverless 7.0.0.

Aplicação	Versão
Apache Spark	3.5.0
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless6.15.0

A tabela a seguir lista as versões das aplicações disponíveis com o EMR Serverless 6.15.0.

Aplicação	Versão
Apache Spark	3.4.1

Aplicação	Versão
Apache Hive	3.1.3
Apache Tez	0.10.2

Notas da versão 6.15.0 do EMR Sem Servidor

- Suporte ao TLS: com as versões 6.15.0 e superiores do Amazon EMR Sem Servidor, habilite a comunicação criptografada por TLS mútua entre trabalhadores em execuções de trabalho do Spark. Quando habilitado, o EMR Sem Servidor gera automaticamente um certificado exclusivo para cada trabalhador, que é provisionado em uma execução de trabalho que os trabalhadores utilizam durante o handshake do TLS para se autenticarem e estabelecerem um canal criptografado para processar dados com segurança. Para obter mais informações sobre a criptografia por TLS mútua, consulte [Criptografia entre trabalhadores](#).

EMR Serverless6.14.0

A tabela a seguir lista as versões das aplicações disponíveis no EMR Serverless 6.14.0.

Aplicação	Versão
Apache Spark	3.4.1
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless6.13.0

A tabela a seguir lista as versões das aplicações disponíveis no EMR Serverless 6.13.0.

Aplicação	Versão
Apache Spark	3.4.1

Aplicação	Versão
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless6.12.0

A tabela a seguir lista as versões das aplicações disponíveis no EMR Serverless 6.12.0.

Aplicação	Versão
Apache Spark	3.4.0
Apache Hive	3.1.3
Apache Tez	0.10.2

EMR Serverless6.11.0

A tabela a seguir lista as versões das aplicações disponíveis no EMR Serverless 6.11.0.

Aplicação	Versão
Apache Spark	3.3.2
Apache Hive	3.1.3
Apache Tez	0.10.2

Notas da versão 6.11.0 do EMR Sem Servidor

- [Acesse recursos do S3 em outras contas](#): com as versões 6.11.0 e posteriores, você pode configurar vários perfis do IAM a serem assumidos ao acessar buckets do Amazon S3 em contas diferentes da AWS no EMR Sem Servidor.

EMR Serverless6.10.0

A tabela a seguir lista as versões das aplicações disponíveis no EMR Serverless 6.10.0.

Aplicação	Versão
Apache Spark	3.3.1
Apache Hive	3.1.3
Apache Tez	0.10.2

Notas da versão 6.10.0 do EMR Sem Servidor

- Em aplicações do EMR Sem Servidor com a versão 6.10.0 ou posterior, o valor padrão da propriedade `spark.dynamicAllocation.maxExecutors` é `infinity`. As versões anteriores são padronizadas para `100`. Para obter mais informações, consulte [Propriedades do trabalho do spark](#).

EMR Serverless6.9.0

A tabela a seguir lista as versões das aplicações disponíveis no EMR Serverless 6.9.0.

Aplicação	Versão
Apache Spark	3.3.0
Apache Hive	3.1.3
Apache Tez	0.10.2

Notas da versão 6.9.0 do EMR Sem Servidor

- A integração do Amazon Redshift para Apache Spark está inclusa nas versões 6.9.0 e posteriores do Amazon EMR. Anteriormente uma ferramenta de código aberto, a integração nativa é um conector do Spark que você pode usar para criar aplicações do Apache Spark que realizam a leitura e a gravação de dados no Amazon Redshift e no Amazon Redshift sem servidor. Para obter

mais informações, consulte [Uso da integração do Amazon Redshift para Apache Spark no Amazon EMR Sem Servidor](#).

- A versão 6.9.0 do EMR Serverless adiciona suporte à arquitetura Graviton2 (arm64). AWS Você pode usar o `architecture` parâmetro para `create-application` e `update-application` APIs para escolher a arquitetura arm64. Para obter mais informações, consulte [Opções de arquitetura do Amazon EMR Sem Servidor](#).
- Agora você pode exportar, importar, consultar e unir tabelas do Amazon DynamoDB diretamente das aplicações do Spark e Hive no EMR Sem Servidor. Para obter mais informações, consulte [Como se conectar ao DynamoDB com o Amazon EMR Sem Servidor](#).

Problemas conhecidos

- Se você usar a integração do Amazon Redshift para Apache Spark e tiver um `time`, `timetz`, `timestamp` ou `timestampz` com precisão de microssegundos no formato Parquet, o conector arredondará os valores de tempo para o valor de milissegundo mais próximo. Como solução alternativa, use o parâmetro `unload_s3_format` do formato de descarregamento de texto.

EMR Serverless6.8.0

A tabela a seguir lista as versões das aplicações disponíveis no EMR Serverless 6.8.0.

Aplicação	Versão
Apache Spark	3.3.0
Apache Hive	3.1.3
Apache Tez	0.9.2

EMR Serverless6.7.0

A tabela a seguir lista as versões das aplicações disponíveis no EMR Serverless 6.7.0.

Aplicação	Versão
Apache Spark	3.2.1

Aplicação	Versão
Apache Hive	3.1.3
Apache Tez	0.9.2

Alterações, melhorias e problemas resolvidos específicos do mecanismo

A tabela a seguir lista um novo recurso específico do mecanismo.

Alteração	Description
Recurso	O programador do Tez agora é compatível com a preempção da tarefa do Tez em vez da preempção do contêiner

EMR Serverless 6.6.0

A tabela a seguir lista as versões das aplicações disponíveis no EMR Serverless 6.6.0.

Aplicação	Versão
Apache Spark	3.2.0
Apache Hive	3.1.2
Apache Tez	0.9.2

Notas da versão inicial do EMR Sem Servidor

- O EMR Sem Servidor é compatível com a classificação de configuração do Spark `spark-defaults`. Essa classificação altera os valores no arquivo XML `spark-defaults.conf` do Spark. As classificações de configuração permitem que você personalize aplicações. Para obter mais informações, consulte [Configurar aplicações](#).
- O EMR Sem Servidor é compatível com as classificações de configuração `hive-site`, `tez-site`, `emrfs-site` e `core-site` do Hive. Essa classificação pode alterar os valores no arquivo

`hive-site.xml` do Hive, no arquivo `tez-site.xml` do Tez, nas configurações do EMRFS do Amazon EMR ou no arquivo `core-site.xml` do Hadoop, respectivamente. As classificações de configuração permitem que você personalize aplicações. Para obter mais informações, consulte [Configurar aplicações](#).

Alterações, melhorias e problemas resolvidos específicos do mecanismo

- A tabela a seguir lista as backports do Hive e do Tez.

Alterações no Hive e no Tez

Alteração	Description
Backport	TEZ-4430 : problema corrigido com a propriedade <code>tez.task.launch.cmd-opts</code>
Backport	HIVE-25971 : foram corrigidos os atrasos no desligamento da tarefa do Tez devido ao grupo de threads em cache aberto

Histórico do documento

A tabela a seguir descreve as alterações importantes na documentação desde a última versão do EMR Sem Servidor. Para obter mais informações sobre as atualizações desta documentação, você pode se tornar assinante de um feed RSS.

Alteração	Descrição	Data
Nova versão prévia	AWS runtime for Apache Spark(pré-visualização do emr-spark-8.0)	21 de novembro de 2025
Nova versão	EMR Serverless7.2.0	25 de julho de 2024
Nova versão	EMR Serverless7.1.0	17 de abril de 2024
Atualize para uma política existente.	Adicionou a nova Sid CloudWatchPolicyStatement e EC2PolicyStatement à EMRServerlessServiceRolePolicy política da Amazon.	25 de janeiro de 2024
Nova versão	EMR Serverless7.0.0	29 de dezembro de 2023
Nova versão	EMR Serverless6.15.0	17 de novembro de 2023
Novo atributo	Configure vários perfis do IAM a serem assumidos ao acessar buckets do Amazon S3 em contas diferentes do EMR Sem Servidor (6.11 e superior)	18 de outubro de 2023
Nova versão	EMR Serverless6.14.0	17 de outubro de 2023

Novo atributo	Configuração padrão de aplicações do EMR Sem Servidor	25 de setembro de 2023
Atualização para as propriedades padrão do Hive	Foram atualizados os valores padrão das propriedades de trabalho do Hive <code>hive.driver.disk</code> , <code>hive.tez.disk.size</code> , <code>hive.tez.auto.reducer.parallelism</code> e <code>tez.grouping.min-size</code> .	12 de setembro de 2023
Nova versão	EMR Serverless6.13.0	11 de setembro de 2023
Nova versão	EMR Serverless6.12.0	21 de julho de 2023
Nova versão	EMR Serverless6.11.0	8 de junho de 2023
Atualizar política de função vinculada ao serviço	Foi atualizado o perfil SLR AmazonEMRServerlessServiceRolePolicy _ para publicar o uso em nível de conta no namespace "AWS/Usage" .	20 de abril de 2023
Disponibilidade geral (GA) do EMR Serverless	Esse é o primeiro lançamento público do EMR Sem Servidor.	1º de junho de 2022

As traduções são geradas por tradução automática. Em caso de conflito entre o conteúdo da tradução e da versão original em inglês, a versão em inglês prevalecerá.