

Pilier Efficacité des performances



Pilier Efficacité des performances: AWS Well-Architected Framework

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

Les marques commerciales et la présentation commerciale d'Amazon ne peuvent pas être utilisées en relation avec un produit ou un service extérieur à Amazon, d'une manière susceptible d'entraîner une confusion chez les clients, ou d'une manière qui dénigre ou discrédite Amazon. Toutes les autres marques commerciales qui ne sont pas la propriété d'Amazon appartiennent à leurs propriétaires respectifs, qui peuvent ou non être affiliés ou connectés à Amazon, ou sponsorisés par Amazon.

Table of Contents

Résumé et introduction	1
Introduction	1
Efficacité des performances	3
Principes de conception	3
Définition	4
Choix d'architecture	5
PERF01-BP01 Découvrez et comprenez les services et fonctionnalités cloud disponibles	5
Directives d'implémentation	6
Ressources	7
PERF01-BP02 Utilisation des recommandations de votre fournisseur de cloud ou d'un partenaire approprié pour en savoir plus sur les modèles d'architecture et les bonnes pratiques	8
Directives d'implémentation	6
Ressources	7
PERF01-BP03 Intégrer les coûts dans les décisions architecturales	10
Directives d'implémentation	6
Ressources	7
PERF01-BP04 Évaluation de l'impact des compromis sur les clients et l'efficacité de l'architecture	12
Directives d'implémentation	6
Ressources	7
PERF01-BP05 Politiques d'utilisation et architectures de référence	14
Directives d'implémentation	6
Ressources	7
PERF01-BP06 Utilisation du benchmarking pour éclairer vos décisions architecturales	16
Directives d'implémentation	6
Ressources	7
PERF01-BP07 Utiliser une approche axée sur les données pour les choix architecturaux	19
Directives d'implémentation	6
Ressources	7
Informatique et matériel	22
PERF02-BP01 Sélectionnez les meilleures options de calcul pour votre charge de travail	22
Directives d'implémentation	6
Étapes d'implémentation	6

Ressources	7
PERF02-BP02 Comprendre la configuration et les fonctionnalités de calcul disponibles	26
Directives d'implémentation	6
Étapes d'implémentation	6
Ressources	7
PERF02-BP03 Collecter des métriques liées au calcul	30
Directives d'implémentation	6
Étapes d'implémentation	6
Ressources	7
PERF02-BP04 Configurer et dimensionner correctement les ressources de calcul	33
Directives d'implémentation	6
Ressources	7
PERF02-BP05 Adaptez dynamiquement vos ressources informatiques	36
Directives d'implémentation	6
Ressources	7
PERF02-BP06 Utiliser des accélérateurs de calcul matériels optimisés	39
Directives d'implémentation	6
Ressources	7
Gestion des données	42
PERF03-BP01 Utilisez un magasin de données spécialement conçu pour répondre au mieux à vos besoins en matière d'accès aux données et de stockage	42
Directives d'implémentation	6
Ressources	7
PERF03-BP02 Évaluer les options de configuration disponibles pour le magasin de données	55
Directives d'implémentation	6
Ressources	7
PERF03-BP03 Collecter et enregistrer les indicateurs de performance du magasin de données	60
Directives d'implémentation	6
Étapes d'implémentation	6
Ressources	7
PERF03-BP04 Mise en œuvre de stratégies pour améliorer les performances des requêtes dans un magasin de données	63
Directives d'implémentation	6
Ressources	7
PERF03-BP05 Implémenter des modèles d'accès aux données qui utilisent la mise en cache ...	65

Directives d'implémentation	6
Ressources	7
Réseau et diffusion de contenu	70
PERF04-BP01 Comprendre l'impact du réseau sur les performances	70
Directives d'implémentation	6
Ressources	7
PERF04-BP02 Évaluer les fonctionnalités réseau disponibles	74
Directives d'implémentation	6
Ressources	7
PERF04-BP03 Choisissez une connectivité dédiée adaptée à votre charge VPN de travail	81
Directives d'implémentation	6
Ressources	7
PERF04-BP04 Utilisez l'équilibrage de charge pour répartir le trafic entre plusieurs ressources	84
Directives d'implémentation	6
Ressources	7
PERF04-BP05 Choisissez les protocoles réseau pour améliorer les performances	89
Directives d'implémentation	6
Ressources	7
PERF04-BP06 Choisissez l'emplacement de votre charge de travail en fonction des exigences du réseau	92
Directives d'implémentation	6
Ressources	7
PERF04-BP07 Optimiser la configuration du réseau en fonction des métriques	98
Directives d'implémentation	6
Ressources	7
Processus et culture	104
PERF05-BP01 Définition d'indicateurs de rendement clés (KPI) pour mesurer l'état et les performances de la charge de travail	106
Directives d'implémentation	6
Étapes d'implémentation	6
Ressources	7
PERF05-BP02 Utilisez des solutions de surveillance pour comprendre les domaines dans lesquels les performances sont les plus critiques	109
Directives d'implémentation	6
Ressources	7

PERF05-BP03 Définir un processus pour améliorer les performances de la charge de travail ..	112
Directives d'implémentation	6
Ressources	7
PERF05-BP04 Testez votre charge de travail	114
Directives d'implémentation	6
Ressources	7
PERF05-BP05 Utiliser l'automatisation pour résoudre de manière proactive les problèmes liés aux performances	116
Directives d'implémentation	6
Ressources	7
PERF05-BP06 Maintenez votre charge de travail et vos services up-to-date	119
Directives d'implémentation	6
Étapes d'implémentation	6
Ressources	7
PERF05-BP07 Vérification des métriques à intervalles réguliers	121
Directives d'implémentation	6
Ressources	7
Conclusion	124
Collaborateurs	125
Suggestions de lecture	126
Révisions du document	127
Avis	129
AWS Glossaire	130

Pilier Efficacité des performances - AWS Well-Architected Framework

Date de publication : 6 novembre 2024 ([Révisions du document](#))

Ce livre blanc porte sur le pilier Efficacité des performances d'AWS Well-Architected Framework. Il fournit des conseils pour aider les clients à appliquer les bonnes pratiques de conception, de distribution et de maintenance des environnements AWS.

Introduction

[AWS Well-Architected Framework](#) vous aide à mesurer le pour et le contre des options qui se présentent lors de la création de charges de travail sur AWS. En utilisant ce cadre, vous apprenez les bonnes pratiques architecturales en matière de conception et d'exploitation de charges de travail fiables, sécurisées, efficaces, économiques et durables dans le cloud. Il vous permet d'évaluer systématiquement vos architectures par rapport aux bonnes pratiques et d'identifier les domaines à améliorer. Nous pensons que le fait d'avoir des charges de travail bien structurées augmente considérablement les chances de réussite métier.

Le cadre repose sur six piliers :

- Excellence opérationnelle
- Sécurité
- Fiabilité
- Efficacité des performances
- Optimisation des coûts
- Durabilité

Ce livre blanc porte sur l'application des principes du pilier Efficacité des performances à vos charges de travail. Dans les environnements sur site traditionnels, il est difficile de bénéficier de performances élevées et durables. En appliquant les principes de ce livre blanc, vous pourrez créer des architectures sur AWS qui fournissent avec efficacité des performances soutenues sur le long terme. Les conseils et les bonnes pratiques présentés dans ce document sont répartis dans cinq domaines clés qui servent de principes directeurs pour la création de solutions cloud performantes sur AWS. Ces domaines d'intérêt sont les suivants :

- [Choix d'architecture](#)
- [Informatique et matériel](#)
- [Gestion des données](#)
- [Réseau et diffusion de contenu](#)
- [Processus et culture](#)

Le présent document est conçu pour ceux et celles qui sont dépositaires de rôles technologiques, comme les directeurs de la technologie, les architectes, les développeurs et les membres de l'équipe d'exploitation. Après avoir lu ce document, vous allez vous familiariser avec les bonnes pratiques et les stratégies d'AWS à utiliser lors de la conception d'architectures cloud performantes.

Efficacité des performances

Le pilier Efficacité des performances englobe la capacité à utiliser efficacement les ressources du cloud pour satisfaire aux exigences système et à maintenir cette efficacité au fur et à mesure que la demande change et que les technologies évoluent.

Rubriques

- [Principes de conception](#)
- [Définition](#)

Principes de conception

Les principes de conception suivants peuvent vous aider à créer des charges de travail efficaces dans le cloud, tout en veillant à ce qu'elles le restent dans la durée.

- Démocratiser les technologies avancées : simplifiez la mise en œuvre de technologies avancées pour votre équipe en déléguant des tâches complexes à votre fournisseur de cloud. Plutôt que de demander à votre équipe informatique de s'informer sur l'hébergement et l'exploitation de nouvelles technologies, envisagez de consommer les technologies en tant que service. Par exemple, l'absence SQL de bases de données, le transcodage multimédia et l'apprentissage automatique sont autant de technologies qui nécessitent une expertise spécialisée. Dans le cloud, ces technologies deviennent des services que votre équipe peut consommer, ce qui lui permet de se consacrer au développement de produits plutôt qu'à l'allocation et à la gestion des ressources.
- Passez à l'international en quelques minutes : le déploiement de votre charge de travail dans plusieurs AWS régions du monde vous permet de réduire le temps de latence et d'offrir une meilleure expérience à vos clients à moindre coût.
- Utilisation d'architectures sans serveur : les architectures sans serveur vous évitent d'exécuter et de gérer des serveurs physiques pour les activités traditionnelles de calcul. Par exemple, les services de stockage sans serveur peuvent agir comme des sites Web statiques (éliminant le besoin de serveurs Web), et les services d'événements peuvent héberger du code. Ainsi, vous supprimez la charge opérationnelle de gestion des serveurs physiques et réduisez les coûts des transactions, car les services gérés fonctionnent à l'échelle du cloud.
- Expérimentation plus fréquente : avec des ressources virtuelles et automatisables, vous pouvez rapidement exécuter des tests comparatifs à l'aide de différents types d'instances, de stockages ou de configurations.

- Envisager la compréhension technique : utilisez l'approche technologique qui correspond le mieux à vos objectifs. Par exemple, tenez compte des modèles d'accès aux données lorsque vous sélectionnez les approches de stockage ou de base de données de votre charge de travail.

Définition

Concentrez-vous sur les domaines suivants pour assurer l'efficacité des performances dans le cloud :

- [Choix d'architecture](#)
- [Informatique et matériel](#)
- [Gestion des données](#)
- [Réseau et diffusion de contenu](#)
- [Processus et culture](#)

Adoptez une approche axée sur les données pour créer une architecture performante. Collectez des données sur tous les aspects de l'architecture, depuis la conception générale jusqu'à la sélection et la configuration des types de ressources.

En revoyant régulièrement vos choix, vous vous assurez de tirer parti de l'évolution constante du AWS Cloud. La surveillance vous offre la garantie d'être informé de tout écart par rapport aux performances attendues. Effectuer des compromis dans votre architecture pour améliorer les performances, comme l'utilisation de la compression, la mise en cache ou l'abaissement des exigences de cohérence.

Choix d'architecture

La solution optimale pour une charge de travail peut varier, et les solutions combinent souvent plusieurs approches. Les charges de travail Well-Architected utilisent plusieurs solutions et permettent d'exploiter différentes fonctionnalités pour améliorer les performances.

De nombreux types et configurations de ressources AWS sont proposés. Il est ainsi plus facile de trouver l'approche qui correspond le mieux à vos besoins. Vous pouvez également rechercher des options qui ne sont pas facilement accessibles avec une infrastructure sur site. Par exemple, un service géré tel qu'Amazon DynamoDB fournit une base de données NoSQL entièrement gérée avec une latence de moins de 10 millisecondes, quelle que soit l'échelle.

Ce domaine d'intérêt partage des conseils et des bonnes pratiques sur la manière de sélectionner des ressources cloud et des modèles d'architecture efficaces et performants.

Bonnes pratiques

- [PERF01-BP01 Découvrez et comprenez les services et fonctionnalités cloud disponibles](#)
- [PERF01-BP02 Utilisation des recommandations de votre fournisseur de cloud ou d'un partenaire approprié pour en savoir plus sur les modèles d'architecture et les bonnes pratiques](#)
- [PERF01-BP03 Intégrer les coûts dans les décisions architecturales](#)
- [PERF01-BP04 Évaluation de l'impact des compromis sur les clients et l'efficacité de l'architecture](#)
- [PERF01-BP05 Politiques d'utilisation et architectures de référence](#)
- [PERF01-BP06 Utilisation du benchmarking pour éclairer vos décisions architecturales](#)
- [PERF01-BP07 Utiliser une approche axée sur les données pour les choix architecturaux](#)

PERF01-BP01 Découvrez et comprenez les services et fonctionnalités cloud disponibles

Découvrez en continu les services et configurations disponibles qui vous aident à prendre de meilleures décisions architecturales et à améliorer l'efficacité des performances de votre architecture de charge de travail.

Anti-modèles courants :

- Vous utilisez le cloud comme centre de données hébergé.

- Vous ne modernisez pas votre application après la migration vers le cloud.
- Vous n'utilisez qu'un seul type de stockage pour tout ce que vous devez conserver.
- Vous utilisez les types d'instances qui correspondent le plus à vos standards actuels. Elles peuvent être de plus grande taille au besoin.
- Vous déployez et gérez les technologies disponibles en tant que services gérés.

Avantages liés au respect de cette bonne pratique : en envisageant de nouveaux services et de nouvelles configurations, vous pourriez être en mesure d'améliorer considérablement vos performances, de réduire les coûts et d'optimiser les efforts requis pour maintenir votre charge de travail. Cela peut également vous aider à accélérer le développement time-to-value des produits compatibles avec le cloud.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

AWS publie en permanence de nouveaux services et fonctionnalités susceptibles d'améliorer les performances et de réduire le coût des charges de travail dans le cloud. Il est essentiel up-to-date de rester fidèle à ces nouveaux services et fonctionnalités pour maintenir l'efficacité des performances dans le cloud. La modernisation de votre architecture de charge de travail vous permet également d'accélérer la productivité, de stimuler l'innovation et de générer de nouvelles opportunités de croissance.

Étapes d'implémentation

- Faites l'inventaire de vos charges de travail logicielles et de l'architecture des services connexes. Déterminez la catégorie de produits sur laquelle vous souhaitez en savoir plus.
- Explorez les AWS offres pour identifier et découvrir les services et options de configuration pertinents qui peuvent vous aider à améliorer les performances et à réduire les coûts et la complexité opérationnelle.
 - [Amazon Web Services Cloud](#)
 - [AWS Académie](#)
 - [Quoi de neuf avec AWS ?](#)
 - [AWS Blog](#)
 - [AWS Générateur de compétences](#)

- [AWS Événements et webinaires](#)
- [AWS Training et certifications](#)
- [AWS Chaîne Youtube](#)
- [AWS Ateliers](#)
- [Communautés AWS](#)
- Utilisez [Amazon Q](#) pour obtenir des informations pertinentes et des conseils sur les services.
- Utilisez des environnements de test (sandbox) (hors production) pour découvrir et tester de nouveaux services sans frais supplémentaires.
- Découvrez en permanence les nouveaux services et fonctionnalités du cloud.

Ressources

Documents connexes :

- [Présentation d'Amazon Web Services](#)
- [EC2Fonctionnalités d'Amazon](#)
- [Apprenez step-by-step avec le plan de formation d'un AWS partenaire](#)
- [AWS Formation et certification](#)
- [Mon parcours d'apprentissage pour devenir architecte de AWS solutions](#)
- [AWS Centre d'architecture](#)
- [AWS Partner Network](#)
- [AWS Bibliothèque de solutions](#)
- [AWS Centre de connaissances](#)
- [Créez des applications modernes sur AWS](#)

Vidéos connexes :

- [AWS re:Invent 2023 - Nouveautés d'Amazon EC2](#)
- [AWS re:Invent 2022 - Réduisez vos coûts d'exploitation et d'infrastructure avec Amazon ECS](#)
- [AWS re:Invent 2023 - Développez avec l'efficacité, l'agilité et l'innovation du cloud avec AWS](#)
- [AWS re:Invent 2022 - Déployez des modèles de machine learning pour l'inférence à des performances élevées et à moindre coût](#)

- [This is my Architecture](#)

Exemples connexes :

- [AWS Exemples](#)
- [AWS SDKExemples](#)

PERF01-BP02 Utilisation des recommandations de votre fournisseur de cloud ou d'un partenaire approprié pour en savoir plus sur les modèles d'architecture et les bonnes pratiques

Utilisez les ressources cloud de l'entreprise, telles que la documentation, les architectes de solutions, les services professionnels ou les partenaires appropriés pour éclairer vos décisions architecturales. Ces ressources vous aident à vérifier et à améliorer votre architecture pour obtenir des performances optimales.

Anti-modèles courants :

- Vous utilisez AWS en tant que fournisseur de cloud ordinaire.
- Vous utilisez les services AWS de manière non conforme à leur utilisation prévue.
- Vous suivez toutes les recommandations sans tenir compte du contexte de votre entreprise.

Avantage de l'établissement de cette bonne pratique : en suivant les recommandations d'un fournisseur de cloud ou d'un partenaire approprié, vous pouvez faire les bons choix architecturaux pour votre charge de travail et vous avez confiance dans vos décisions.

Niveau de risque encouru si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

AWS propose un large éventail de recommandations, documentations et ressources qui peuvent vous aider à générer et à gérer des charges de travail cloud efficaces. La documentation AWS fournit des exemples de code, des tutoriels et des explications détaillées sur les services. Outre la documentation, AWS propose des programmes de formation et de certification, des architectes de solutions et des services professionnels qui peuvent aider les clients à explorer différents aspects des services cloud et à mettre en œuvre une architecture cloud efficace sur AWS.

Tirez parti de ces ressources pour obtenir des informations précieuses et des bonnes pratiques, gagner du temps et obtenir de meilleurs résultats dans le AWS Cloud.

Étapes d'implémentation

- Consultez la documentation et les recommandations AWS et suivez les bonnes pratiques. Ces ressources peuvent vous aider à choisir et à configurer efficacement les services, ainsi qu'à améliorer les performances.
 - [Documentation AWS](#) (comme les guides d'utilisation et les livres blancs)
 - [Blog AWS](#)
 - [AWS Training et certifications](#)
 - [Chaîne YouTube AWS](#)
- Participez à des événements partenaires AWS (tels que les sommets mondiaux AWS, les groupes d'utilisateurs, re:Invent AWS et les ateliers) pour découvrir les bonnes pratiques d'utilisation des services AWS auprès des experts AWS.
 - [Apprentissage étape par étape grâce à un plan de formation pour les partenaires AWS](#)
 - [Événements et webinaires AWS](#)
 - [Ateliers AWS](#)
 - [Communautés AWS](#)
- Contactez AWS pour obtenir de l'aide lorsque vous avez besoin de conseils ou d'informations supplémentaires sur le produit. AWS Les architectes de solutions et les [services professionnels AWS](#) prodiguent des conseils pour la mise en œuvre de solutions. [AWS Les partenaires](#) apportent une expertise AWS pour vous aider à gagner en agilité et favoriser l'innovation au sein de votre entreprise.
- Utilisez [Support](#) si vous avez besoin d'une assistance technique pour utiliser un service de manière efficace. [Nos plans de support](#) sont conçus pour vous fournir la bonne combinaison d'outils et l'accès à une expertise afin que vous puissiez réussir avec AWS tout en optimisant les performances, en gérant les risques et en maîtrisant les coûts.

Ressources

Documents connexes :

- [Centre d'architecture AWS](#)
- [AWS Partner Network](#)

- [Bibliothèque de solutions AWS](#)
- [Centre de connaissances AWS](#)
- [AWS Enterprise Support](#)

Vidéos connexes :

- [This is my Architecture](#)
- [AWS re:Invent 2023 - Advanced event-driven patterns with Amazon EventBridge](#)
- [AWS re:Invent 2023 - Implementing distributed design patterns on AWS](#)
- [AWS re:Invent 2023 - Application architecture as code](#)

Exemples connexes :

- [Exemples AWS](#)
- [Exemples de kit SDK AWS](#)
- [Architecture de référence pour l'analytique AWS](#)

PERF01-BP03 Intégrer les coûts dans les décisions architecturales

Tenez compte des coûts dans vos décisions architecturales afin d'améliorer l'utilisation des ressources et l'efficacité des performances de votre charge de travail cloud. Lorsque vous êtes conscient des implications financières de votre charge de travail cloud, vous êtes plus susceptible de tirer parti de ressources efficaces et de réduire les pratiques inutiles.

Anti-modèles courants :

- Vous n'utilisez qu'une seule famille d'instances.
- Vous n'évaluez pas les solutions sous licence par rapport aux solutions open source.
- Vous ne définissez pas de stratégies de cycle de vie pour le stockage.
- Vous ne passez pas en revue les nouveaux services et fonctionnalités du AWS Cloud.
- Vous utilisez uniquement le stockage par blocs.

Avantages liés au respect de cette bonne pratique : en tenant compte des coûts dans vos prises de décision, vous pouvez utiliser des ressources plus efficaces et explorer d'autres investissements.

Niveau de risque encouru si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

L'optimisation des charges de travail en matière de coûts peut améliorer l'utilisation des ressources et éviter le gaspillage dans une charge de travail cloud. La prise en compte des coûts dans les décisions architecturales implique généralement de dimensionner correctement les composants de la charge de travail et de renforcer l'élasticité, ce qui se traduit par une amélioration de l'efficacité des performances de la charge de travail cloud.

Étapes d'implémentation

- Fixez des objectifs de coûts tels que des limites budgétaires pour votre charge de travail cloud.
- Identifiez les composants clés (tels que les instances et le stockage) qui augmentent le coût de votre charge de travail. [Calculateur de tarification AWS](#) et [AWS Cost Explorer](#) vous permettent d'identifier les principaux facteurs de coûts dans votre charge de travail.
- Comprenez les [modèles de tarification](#) dans le cloud, tels que la demande, les instances réservées, les Savings Plans et les instances ponctuelles.
- Utilisez les [bonnes pratiques d'optimisation des coûts de Well-Architected](#) pour optimiser ces composants clés en matière de coûts.
- Surveillez et analysez en permanence les coûts afin d'identifier les opportunités d'optimisation des coûts dans votre charge de travail.
 - Utilisez les [budgets AWS](#) pour recevoir des alertes en cas de coûts inacceptables.
 - Utilisez [AWS Compute Optimizer](#) ou [AWS Trusted Advisor](#) pour obtenir des recommandations en matière d'optimisation des coûts.
 - Utilisez la [détection des anomalies de coûts AWS](#) pour obtenir une détection automatisée des anomalies de coûts et une analyse des causes profondes.

Ressources

Documents connexes :

- [Qu'est-ce que AWS Billing and Cost Management ?](#)
- [Optimisation des coûts avec AWS](#)
- [Choix d'une stratégie de gestion des AWS coûts](#)
- [Guide de gestion des AWS coûts pour débutants](#)

- [Présentation détaillée du tableau de bord Cost Intelligence Dashboard](#)
- [Centre d'architecture AWS](#)
- [Bibliothèque de solutions AWS](#)
- [Centre de connaissances AWS](#)

Vidéos connexes :

- [This is my Architecture](#)
- [AWS re:Invent 2023 - Nouveautés en matière d'optimisation des coûts AWS](#)
- [AWS re:Invent 2023 - Optimisez les coûts et les performances et suivez les progrès en matière d'atténuation](#)
- [AWS re:Invent 2023 - meilleures pratiques en matière d'optimisation des coûts AWS de stockage](#)
- [AWS re:Invent 2023 - Optimisez les coûts dans vos environnements multi-comptes](#)

Exemples connexes :

- [AWS Compute Optimizer Code de démonstration](#)
- [Atelier d'optimisation des coûts](#)
- [Playbooks de mise en œuvre technique de la gestion financière dans le cloud](#)
- [Optimisation du démarrage : ajustement des performances des applications pour une efficacité maximale](#)
- [Atelier d'optimisation sans serveur \(performances et coûts\)](#)
- [Mise à l'échelle d'architectures rentables](#)

PERF01-BP04 Évaluation de l'impact des compromis sur les clients et l'efficacité de l'architecture

Lors de l'évaluation des améliorations liées à la performance, identifiez les choix qui affectent vos clients et l'efficacité de la charge de travail. Par exemple, si l'utilisation d'un magasin de données clé-valeur augmente les performances du système, il est important d'évaluer l'impact de la nature constante de cette modification à terme sur les clients.

Anti-modèles courants :

- Vous supposez que tous les gains de performances doivent être mis en œuvre, même s'il existe des compromis en termes d'implémentation.
- Vous n'évaluez les modifications apportées aux charges de travail que lorsqu'un problème de performances a atteint un point critique.

Avantages liés au respect de cette bonne pratique : lorsque vous évaluez les améliorations potentielles liées aux performances, vous devez décider si les compromis concernant les modifications sont compatibles avec les exigences de charge de travail. Dans certains cas, vous devrez peut-être mettre en place des contrôles supplémentaires pour compenser les compromis.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

Identifiez les domaines critiques de votre architecture en termes de performances et d'impact sur les clients. Déterminez la façon dont vous pouvez apporter des améliorations ainsi que les compromis que ces améliorations entraînent et la façon dont ils affectent le système et l'expérience de l'utilisateur. Par exemple, la mise en œuvre de la mise en cache des données permet d'améliorer de manière significative les performances, mais nécessite une stratégie précise concernant la manière et le moment où mettre à jour ou invalider les données mises en cache pour empêcher un comportement incorrect du système.

Étapes d'implémentation

- Comprenez vos exigences en matière de charge de travail et vos SLA.
- Définissez clairement les facteurs d'évaluation. Les facteurs peuvent être liés au coût, à la fiabilité, à la sécurité et aux performances de votre charge de travail.
- Sélectionnez l'architecture et les services qui répondent à vos besoins.
- Menez des expériences et des démonstrations de faisabilité (POC) afin d'évaluer les facteurs de compromis et l'impact sur les clients et l'efficacité de l'architecture. En général, les charges de travail hautement disponibles, performantes et sécurisées consomment davantage de ressources cloud tout en offrant une meilleure expérience client. Comprenez les compromis entre la complexité, les performances et les coûts de votre charge de travail. Généralement, la priorisation de deux des facteurs se fait au détriment du troisième.

Ressources

Documents connexes :

- [Bibliothèque Amazon Builders' Library](#)
- [KPI QuickSight](#)
- [Amazon CloudWatch RUM](#)
- [Documentation X-Ray](#)
- [Comprenez les modèles de résilience et les compromis pour concevoir une architecture efficace dans le cloud](#)

Vidéos connexes :

- [Optimize applications through via Amazon CloudWatch RUM](#)
- [AWSre:Invent 2023 - Capacity, availability, cost efficiency: Pick three](#)
- [AWS re:Invent 2023 - Advanced integration patterns & trade-offs for loosely coupled systems](#)

Exemples connexes :

- [Measure page load time with Amazon CloudWatch Synthetics](#)
- [Client Web Amazon CloudWatch RUM](#)

PERF01-BP05 Politiques d'utilisation et architectures de référence

Utilisez les stratégies internes et les architectures de référence existantes lors de la sélection des services et des configurations en vue d'augmenter votre efficacité lorsque vous concevez et mettez en œuvre votre charge de travail.

Anti-modèles courants :

- Vous autorisez un large éventail de technologies qui peuvent avoir un impact sur les frais généraux de gestion de votre entreprise.

Avantages liés au respect de cette bonne pratique : l'établissement d'une stratégie pour les choix d'architecture, de technologie et de fournisseur permet de prendre des décisions rapidement.

Niveau de risque encouru si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Le fait de disposer de stratégies internes en matière de sélection des ressources et de l'architecture fournit des normes et des directives à suivre lors des choix architecturaux. Ces directives simplifient le processus de prise de décision lors du choix du bon service cloud et peuvent contribuer à améliorer l'efficacité des performances. Déployez votre charge de travail à l'aide de stratégies ou d'architectures de référence. Intégrez les services à votre déploiement dans le cloud. Utilisez ensuite vos tests de performance pour vérifier que vous pouvez continuer à répondre à vos exigences de performance.

Étapes d'implémentation

- Comprenez clairement les exigences de votre charge de travail cloud.
- Passez en revue les stratégies internes et externes pour identifier les plus pertinentes.
- Utilisez les architectures de référence appropriées fournies par AWS ou les bonnes pratiques de votre secteur.
- Créez un continuum composé de stratégies, de normes, d'architectures de référence et de directives normatives pour les situations courantes. Vos équipes pourront ainsi agir plus rapidement. Adaptez les ressources à votre secteur d'activité, le cas échéant.
- Validez ces stratégies et architectures de référence pour votre charge de travail dans les environnements de test (sandbox).
- up-to-dateRespectez les normes et les AWS mises à jour du secteur pour vous assurer que vos politiques et architectures de référence contribuent à optimiser votre charge de travail dans le cloud.

Ressources

Documents connexes :

- [Centre d'architecture AWS](#)
- [AWS Partner Network](#)
- [Bibliothèque de solutions AWS](#)
- [Centre de connaissances AWS](#)

- [AWS Blogue d'architecture](#)

Vidéos connexes :

- [This is my Architecture](#)
- [AWS re:Invent 2022 - Accélérez la création de valeur pour votre entreprise grâce à une architecture de SAP référence AWS](#)

Exemples connexes :

- [Exemples AWS](#)
- [AWS SDKExemples](#)

PERF01-BP06 Utilisation du benchmarking pour éclairer vos décisions architecturales

Définissez des points de référence pour les performances d'une charge de travail existante afin de comprendre ses performances sur le cloud et prendre des décisions architecturales sur la base de ces données.

Anti-modèles courants :

- Vous comptez sur des points de référence courants qui ne reflètent pas les caractéristiques de votre charge de travail.
- Vous utilisez les commentaires et la perception des clients comme seule référence.

Avantages de l'établissement de cette bonne pratique : le benchmarking de votre implémentation actuelle vous permet de mesurer les améliorations de performance.

Niveau de risque encouru si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Utilisez la définition de points de référence avec des tests synthétiques pour évaluer les performances des composants de votre charge de travail. Le benchmarking est généralement plus rapide à configurer que les tests de charge. Il est utilisé pour évaluer la technologie pour un

composant en particulier. Le benchmarking est souvent utilisé au début d'un nouveau projet, lorsque vous n'avez pas de solution complète pour le test de charge.

Vous pouvez créer vos propres tests de performances, ou bien utiliser un test conforme aux normes du secteur, comme le [TPC-DS](#) pour comparer vos charges de travail. Les points de référence du secteur sont utiles lorsque vous comparez différents environnements. Les points de référence personnalisés sont utiles pour cibler certains types d'opérations que vous souhaitez effectuer dans votre architecture.

Avec le benchmarking, il est important de préparer votre environnement de test pour obtenir des résultats valides. Exécutez plusieurs fois le même point de référence pour vous assurer d'avoir capturé toute variabilité au fil du temps.

Étant donné que les points de référence sont généralement plus rapides à exécuter que les tests de charge, ils peuvent être utilisés plus tôt dans le pipeline de déploiement et fournir un retour rapide sur les écarts de performances. Lorsque vous évaluez un changement important dans un composant ou un service, un point de référence peut être un moyen rapide pour voir si la modification a un intérêt. L'utilisation du benchmarking avec un test de charge est essentielle, car un test de charge vous indique comment votre charge de travail se comporte dans un environnement de production.

Étapes d'implémentation

- Planification et définition :
 - Définissez les objectifs, la base de référence, les scénarios de test, les métriques (telles que l'utilisation du processeur, la latence ou le débit) et les indicateurs de rendement clés de votre test de performances.
 - Concentrez-vous sur les exigences des utilisateurs en matière d'expérience utilisateur et sur des facteurs tels que le temps de réponse et l'accessibilité.
 - Identifiez un outil de benchmarking adapté à votre charge de travail. Vous pouvez utiliser des services AWS tels qu'[Amazon CloudWatch](#) ou un outil tiers compatible avec votre charge de travail.
- Configuration et instrumentation :
 - Configurez votre environnement et vos ressources.
 - Mettez en œuvre la surveillance et la journalisation pour capturer les résultats des tests.
- Comparaison et surveillance :
 - Effectuez vos tests de performances et surveillez les métriques pendant le test.
- Analyse et documentation :

- Documentez votre processus de benchmarking et vos résultats.
- Analysez les résultats pour identifier les goulots d'étranglement, les tendances et les domaines d'amélioration.
- Utilisez les résultats des tests pour prendre des décisions architecturales et ajuster votre charge de travail. Cet ajustement peut impliquer la modification des services ou l'adoption de nouvelles fonctionnalités.
- Optimisation et répétition :
 - Ajustez les configurations et les allocations des ressources en fonction de vos critères de référence.
 - Testez à nouveau votre charge de travail après ajustement pour valider vos améliorations.
 - Documentez vos conclusions et répétez le processus pour identifier d'autres domaines d'amélioration.

Ressources

Documents connexes :

- [Centre d'architecture AWS](#)
- [AWS Partner Network](#)
- [Bibliothèque de solutions AWS](#)
- [Centre de connaissances AWS](#)
- [Amazon CloudWatch RUM](#)
- [Amazon CloudWatch Synthetics](#)
- [Flux de travail génomiques, partie 5 : benchmarking automatisé](#)
- [Évaluation et optimisation du déploiement des points de terminaison dans Amazon SageMaker AI JumpStart](#)

Vidéos connexes :

- [AWS re:Invent 2023 - Benchmarking AWS Lambda cold starts](#)
- [Benchmarking stateful services in the cloud](#)
- [This is my Architecture](#)
- [Optimize applications through via Amazon CloudWatch RUM](#)

- [Présentation d'Amazon CloudWatch Synthetics](#)

Exemples connexes :

- [Exemples AWS](#)
- [Exemples de kit SDK AWS](#)
- [Tests de charge distribuée](#)
- [Mesure du temps de chargement des pages avec Amazon CloudWatch Synthetics](#)
- [Client Web Amazon CloudWatch RUM](#)

PERF01-BP07 Utiliser une approche axée sur les données pour les choix architecturaux

Définissez une approche orientée données claire pour les choix architecturaux afin de vérifier que les services et configurations cloud appropriés sont utilisés pour répondre aux besoins spécifiques de votre entreprise.

Anti-modèles courants :

- Vous supposez que votre architecture actuelle est statique et ne devrait pas être mise à jour au fil du temps.
- Vos choix architecturaux sont basés sur des suppositions et des hypothèses.
- Vous introduisez des modifications d'architecture au fil du temps sans justification.

Avantages liés au respect de cette bonne pratique : en adoptant une approche bien définie pour les choix architecturaux, vous utilisez les données pour influencer la conception de votre charge de travail et prendre des décisions éclairées au fil du temps.

Niveau de risque encouru si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Mobilisez l'expérience et l'expertise des ressources cloud internes ou faites appel à des ressources externes, comme des cas d'utilisation publiés ou des livres blancs pour choisir les ressources et services dans votre architecture. Vous devriez disposer d'un processus bien défini qui encourage

l'expérimentation et le benchmarking avec les services qui pourraient être utilisés dans votre charge de travail.

Les backlogs relatifs aux charges de travail critiques doivent non seulement comprendre des témoignages d'utilisateurs proposant des fonctionnalités pertinentes pour les entreprises et les utilisateurs, mais également des récits techniques qui constituent une piste architecturale pour la charge de travail. Cette piste s'inspire des nouvelles avancées technologiques et des nouveaux services et les adopte sur la base de données et de justifications appropriées. Cela permet de vérifier que l'architecture reste pérenne et ne stagne pas.

Étapes d'implémentation

- Collaborez avec les principales parties prenantes pour définir les exigences en matière de charge de travail, y compris les considérations relatives aux performances, à la disponibilité et aux coûts. Tenez compte de facteurs tels que le nombre d'utilisateurs et le modèle d'utilisation de votre charge de travail.
- Créez une piste architecturale ou un backlog technologique qui est axé en priorité sur le backlog fonctionnel.
- Évaluez les différents services cloud (pour en savoir plus, consultez [PERF01-BP01 Découvrez et comprenez les services et fonctionnalités cloud disponibles](#)).
- Explorez les différents modèles architecturaux, tels que les microservices ou le modèle sans serveur, qui répondent à vos exigences en termes de performances (pour en savoir plus, consultez [PERF01-BP02 Utilisation des recommandations de votre fournisseur de cloud ou d'un partenaire approprié pour en savoir plus sur les modèles d'architecture et les bonnes pratiques](#)).
- Consultez d'autres équipes, des diagrammes d'architecture et des ressources, telles que AWS Solutions Architects, [AWS Architecture Center](#), etc. [AWS Partner Network](#), pour vous aider à choisir l'architecture adaptée à votre charge de travail.
- Définissez des métriques de performances telles que le débit et le temps de réponse qui peuvent vous aider à évaluer les performances de votre charge de travail.
- Testez et utilisez des métriques définies pour valider les performances de l'architecture sélectionnée.
- Surveillez en continu les performances et effectuez les ajustements nécessaires pour maintenir un niveau optimal de performance pour votre architecture.

- Documentez l'architecture que vous avez sélectionnée et les décisions que vous avez prises comme référence pour les futures mises à jour et les futurs apprentissages.
- Vérifiez en permanence l'approche de sélection de l'architecture et mettez-la à jour en fonction des apprentissages, des nouvelles technologies et des métriques indiquant un changement nécessaire ou un problème dans l'approche actuelle.

Ressources

Documents connexes :

- [Bibliothèque de solutions AWS](#)
- [Centre de connaissances AWS](#)
- [Modèles architecturaux sur lesquels créer End-to-End des applications basées sur les données AWS](#)

Vidéos connexes :

- [This is my Architecture](#)
- [AWS re:Invent 2021 - L'entreprise axée sur les données : passer de la vision à la valeur](#)
- [AWS re:Invent 2022 - Fournir des architectures durables et performantes](#)
- [AWS re:Invent 2023 - Optimisez les coûts et les performances et suivez les progrès en matière d'atténuation](#)
- [AWS re:Invent 2022 - AWS optimisation : étapes réalisables pour des résultats immédiats](#)

Exemples connexes :

- [Exemples AWS](#)
- [AWS SDKExemples](#)

Informatique et matériel

Le choix d'une solution de calcul optimale pour une charge de travail particulière peut varier selon la conception de l'application, les modèles d'utilisation et les paramètres de configuration. Les architectures peuvent utiliser différentes solutions de calcul pour divers composants et permettent différentes fonctionnalités pour améliorer les performances. Le choix d'une solution de calcul inadaptée à une architecture peut nuire à ses performances.

Ce domaine d'intérêt partage des conseils et de bonnes pratiques sur la manière d'identifier et d'optimiser les options de calcul pour l'efficacité des performances dans le cloud.

Bonnes pratiques

- [PERF02-BP01 Sélectionnez les meilleures options de calcul pour votre charge de travail](#)
- [PERF02-BP02 Comprendre la configuration et les fonctionnalités de calcul disponibles](#)
- [PERF02-BP03 Collecter des métriques liées au calcul](#)
- [PERF02-BP04 Configurer et dimensionner correctement les ressources de calcul](#)
- [PERF02-BP05 Adaptez dynamiquement vos ressources informatiques](#)
- [PERF02-BP06 Utiliser des accélérateurs de calcul matériels optimisés](#)

PERF02-BP01 Sélectionnez les meilleures options de calcul pour votre charge de travail

La sélection de l'option de calcul la mieux adaptée à votre charge de travail vous permet d'améliorer les performances, de réduire les coûts d'infrastructure inutiles et de diminuer les efforts opérationnels nécessaires pour maintenir votre charge de travail.

Anti-modèles courants :

- Vous utilisez la même option de calcul que celle utilisée sur site.
- Vous manquez de connaissances sur les options, les fonctionnalités et les solutions de cloud computing et sur la manière dont elles pourraient améliorer vos performances de calcul.
- Vous surprovisionnez une option de calcul existante pour répondre aux exigences de mise à l'échelle ou de performances, alors qu'une autre option de calcul s'alignerait plus précisément sur les caractéristiques de votre charge de travail.

Avantages liés au respect de cette bonne pratique : en identifiant les exigences de calcul et en les comparant aux options disponibles, vous pouvez optimiser votre charge de travail en termes de ressources.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

Pour optimiser vos charges de travail dans le cloud en termes d'efficacité des performances, il est important de sélectionner les options de calcul les plus adaptées à votre cas d'utilisation et à vos exigences de performance. AWS fournit une variété d'options de calcul adaptées aux différentes charges de travail dans le cloud. Par exemple, vous pouvez utiliser [Amazon EC2](#) pour lancer et gérer des serveurs virtuels, [AWS Lambda](#) pour exécuter du code sans avoir à approvisionner ou à gérer des serveurs, [Amazon ECS](#) ou [Amazon EKS](#) pour exécuter et gérer des conteneurs, ou [AWS Batch](#) pour traiter de gros volumes de données en parallèle. En fonction de vos besoins en termes de mise à l'échelle et de calcul, vous devez choisir et configurer la solution de calcul optimale pour votre situation. Vous pouvez également envisager d'utiliser plusieurs types de solutions de calcul dans une seule charge de travail, car chacune présente ses avantages et ses inconvénients.

Les étapes suivantes vous guident dans la sélection des options de calcul adaptées aux caractéristiques de votre charge de travail et à vos exigences de performances.

Étapes d'implémentation

- Comprenez les exigences de calcul de votre charge de travail. Les exigences clés à prendre en compte incluent les besoins de traitement, les modèles de trafic, les modèles d'accès aux données, les besoins de mise à l'échelle et les exigences de latence.
- Découvrez les différents [services AWS informatiques](#) adaptés à votre charge de travail. Pour de plus amples informations, veuillez consulter [PERF01-BP01 Découvrez et comprenez les services et fonctionnalités cloud disponibles](#). Voici quelques options de calcul AWS clés, leurs caractéristiques et leurs cas d'utilisation courants :

AWS service	Principales caractéristiques	Cas d'utilisation courants
Amazon Elastic Compute Cloud (AmazonEC2)	Possède une option dédiée pour le matériel, les exigences de licence, une large sélection de différent	Migration « lift-and-shift », application monolithique, environnements hybrides, applications d'entreprise

AWS service	Principales caractéristiques	Cas d'utilisation courants
	es familles d'instances, les types de processeurs et les accélérateurs de calcul	
Amazon Elastic Container Service (AmazonECS) , Amazon Elastic Kubernetes Service (Amazon) EKS	Déploiement facile, environnements cohérents, évolutivité	Microservices, environnements hybrides
AWS Lambda	Service de calcul sans serveur qui exécute du code en réponse à des événements et gère automatiquement les ressources de calcul sous-jacentes.	Microservices, applications basées sur les événements
AWS Batch	Provisionne et fait évoluer de manière efficace et dynamique Amazon Elastic Container Service (AmazonECS) , Amazon Elastic Kubernetes Service (EKSAmazon) et les ressources de calcul AWS Fargate , avec la possibilité d'utiliser des instances à la demande ou ponctuelles en fonction des exigences de votre poste	HPC, train les modèles ML
Amazon Lightsail	Application Linux et Windows préconfigurée pour exécuter de petites charges de travail	Applications web simples, site web personnalisé

- Évaluez les coûts (tels que le tarif horaire ou le transfert de données) et les frais de gestion (tels que l'application de correctifs et la mise à l'échelle) associés à chaque option de calcul.

- Réalisez des tests et procédez au benchmarking dans un environnement hors production afin d'identifier l'option de calcul la mieux adaptée à vos exigences en termes de charge de travail.
- Après avoir testé et identifié votre nouvelle solution de calcul, planifiez votre migration et validez vos métriques de performance.
- Utilisez AWS des outils de surveillance tels qu'[Amazon CloudWatch](#) et des services d'optimisation [AWS Compute Optimizer](#) pour optimiser en permanence vos ressources informatiques en fonction de modèles d'utilisation réels.

Ressources

Documents connexes :

- [Cloud Compute with AWS](#)
- [Types d'EC2instances Amazon](#)
- [EKSConteneurs Amazon : Amazon EKS Worker Nodes](#)
- [Amazon ECS Containers : instances de ECS conteneurs Amazon](#)
- [Fonctions : configuration des fonctions Lambda](#)
- [Recommandations pour les conteneurs](#)
- [Recommandations pour les modèles sans serveur](#)

Vidéos connexes :

- [AWS re:Invent 2023 - AWS Graviton : le meilleur rapport qualité/prix pour vos charges de travail AWS](#)
- [AWS re:Invent 2023 - Nouvelles fonctionnalités d'IA générative d'Amazon Elastic Compute Cloud dans AMS](#)
- [AWS re:Invent 2023 - What's new with Amazon Elastic Compute Cloud](#)
- [AWS re:Invent 2023 - Smart savings: Amazon Elastic Compute Cloud cost-optimization strategies](#)
- [AWS re:Invent 2021 - Powering next-gen Amazon Elastic Compute Cloud: Deep dive on the Nitro System](#)
- [AWS re:Invent 2019 - Optimisez les performances et les coûts de votre calcul AWS](#)
- [AWS re:Invent 2019 - Amazon Elastic Compute Cloud foundations](#)

- [AWS re:Invent 2022 - Déployez des modèles ML pour l'inférence à des performances élevées et à moindre coût](#)
- [AWS re:Invent 2019 - Optimisez les performances et les coûts de votre calcul AWS](#)
- [EC2Fondations Amazon](#)
- [Deploy ML models for inference at high performance and low cost](#)

Exemples connexes :

- [Migration de l'application Web vers des conteneurs](#)
- [Exécution d'un modèle Hello World sans serveur](#)
- [EKSAtelier Amazon](#)
- [EC2Atelier Amazon](#)
- [Charges de travail efficaces et résilientes avec l'autoscaling Amazon EC2 Auto Scaling](#)
- [Migrer vers AWS Graviton avec Container Services](#)

PERF02-BP02 Comprendre la configuration et les fonctionnalités de calcul disponibles

Découvrez les options et les fonctionnalités de configuration disponibles pour votre service de calcul qui vous aideront à allouer la quantité de ressources appropriée et à améliorer l'efficacité des performances.

Anti-modèles courants :

- Vous ne comparez pas les options de calcul ni les familles d'instances disponibles avec les caractéristiques de la charge de travail.
- Vous surprovisionnez les ressources de calcul pour répondre aux pics de demande.

Avantages de l'établissement de cette meilleure pratique : familiarisez-vous avec les fonctionnalités et les configurations de AWS calcul afin de pouvoir utiliser une solution informatique optimisée pour répondre aux caractéristiques et aux besoins de votre charge de travail.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Chaque solution de calcul dispose de configurations et de fonctionnalités uniques pour prendre en charge différentes caractéristiques et exigences de charge de travail. Découvrez comment ces options soutiennent votre charge de travail et déterminez celles qui sont optimales pour votre système. Ces options incluent par exemple la famille d'instances, les tailles, les fonctionnalités (E/S)GPU, le bursting, les délais d'expiration, la taille des fonctions, les instances de conteneur et la simultanéité. Si votre charge de travail utilise la même option de calcul depuis plus de quatre semaines et que vous pensez que les caractéristiques resteront les mêmes à l'avenir, vous pouvez vérifier si votre option de calcul actuelle est adaptée aux charges de travail CPU et du point de vue de la mémoire. [AWS Compute Optimizer](#)

Étapes d'implémentation

- Comprenez les exigences en matière de charge de travail (comme les CPU besoins, la mémoire et le temps de latence).
- Consultez AWS la documentation et les meilleures pratiques pour découvrir les options de configuration recommandées qui peuvent contribuer à améliorer les performances de calcul. Voici quelques options de configuration clés à prendre en compte :

Option de configuration	Exemples
Type d'instance	<ul style="list-style-type: none">• Les instances optimisées pour le calcul sont idéales pour les charges de travail qui nécessitent un rapport v/mémoire CPU élevé.• Les instances à mémoire optimisée offrent de grandes quantités de mémoire pour soutenir les charges de travail gourmandes en mémoire.• Les instances optimisées pour le stockage sont conçues pour les charges de travail qui nécessitent un accès séquentiel élevé en lecture et en écriture (IOPS) au stockage local.

Option de configuration	Exemples
Modèle de tarification	<ul style="list-style-type: none">• Les instances à la demande vous permettent d'utiliser la capacité de calcul à l'heure ou à la seconde, sans engagement à long terme. Ces instances sont idéales pour dépasser les besoins de base en matière de performances.• Les Savings Plans permettent de réaliser des économies importantes par rapport aux instances à la demande, en échange d'un engagement à utiliser une quantité spécifique de puissance de calcul pour une période d'un ou de trois ans.• Les instances Spot vous permettent de tirer parti de la capacité d'instance inutilisée à un prix réduit pour vos charges de travail sans état et tolérantes aux pannes.
Auto Scaling	Utilisez la configuration Auto Scaling pour faire correspondre les ressources de calcul aux modèles de trafic.
Dimensionnement	<ul style="list-style-type: none">• Utilisez Compute Optimizer pour obtenir des recommandations basées sur le machine learning sur la configuration de calcul qui correspond le mieux à vos caractéristiques de calcul.• Utilisez AWS Lambda Power Tuning pour sélectionner la meilleure configuration pour votre fonction Lambda.

Option de configuration	Exemples
Accélérateurs de calcul matériels	<ul style="list-style-type: none">• Les instances de calcul accéléré exécutent des fonctions telles que le traitement graphique ou la mise en correspondance de modèles de données de manière plus efficace que les alternatives CPU basées sur des données.• Pour les charges de travail liées à l'apprentissage automatique, profitez d'un matériel spécialement conçu pour votre charge de travail, tel que AWS Trainium, Inferentia et Amazon AWS EC2 DL1

Ressources

Documents connexes :

- [Cloud Compute with AWS](#)
- [Types d'EC2 instances Amazon](#)
- [Contrôle de l'état du processeur pour votre EC2 instance Amazon](#)
- [EKSConteneurs Amazon : Amazon EKS Worker Nodes](#)
- [Amazon ECS Containers : instances de ECS conteneurs Amazon](#)
- [Fonctions : configuration des fonctions Lambda](#)

Vidéos connexes :

- [AWS re:Invent 2023 — AWS Graviton : le meilleur rapport qualité/prix pour vos charges de travail AWS](#)
- [AWS re:Invent 2023 — Nouvelles fonctionnalités d'IA EC2 générative d'Amazon dans AWS Management Console](#)
- [AWS re:Invent 2023 — Nouveautés d'Amazon EC2](#)
- [AWS re:Invent 2023 — Économies intelligentes : stratégies d'optimisation des coûts d'Amazon EC2](#)

- [AWS re:Invent 2021 — Au service d'EC2 Amazon de nouvelle génération : étude approfondie du système Nitro](#)
- [AWS re:Invent 2019 — Amazon Foundations EC2](#)
- [AWS re:Invent 2022 — Optimisation des performances et EKS des coûts d'Amazon AWS](#)

Exemples connexes :

- [Code de démonstration de Compute Optimizer](#)
- [Atelier sur les instances Amazon EC2 Spot](#)
- [Charges de travail efficaces et résilientes avec Amazon EC2 AWS Auto Scaling](#)
- [Atelier pour développeurs Graviton](#)
- [AWS journée d'immersion pour les charges de travail Microsoft](#)
- [AWS journée d'immersion pour les charges de travail Linux](#)
- [AWS Compute Optimizer Code de démonstration](#)
- [EKSAtelier Amazon](#)

PERF02-BP03 Collecter des métriques liées au calcul

Enregistrez et suivez les métriques liées au calcul pour mieux comprendre comment fonctionnent vos ressources de calcul et améliorer leurs performances et leur utilisation.

Anti-modèles courants :

- Vous utilisez uniquement la recherche manuelle des fichiers journaux pour les métriques.
- Vous n'utilisez que les métriques par défaut enregistrées par votre logiciel de surveillance.
- Vous n'examinez les métriques qu'en cas de problème.

Avantages liés au respect de cette bonne pratique : en collectant des métriques liées aux performances, vous pouvez aligner les performances des applications sur les exigences de l'entreprise afin de garantir que vous répondez à vos besoins en matière de charge de travail. Cela peut également vous aider à améliorer en continu les performances et l'utilisation des ressources de votre charge de travail.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

Les charges de travail cloud peuvent générer de gros volumes de données telles que des métriques, des journaux et des événements. Dans ce contexte AWS Cloud, la collecte de métriques est une étape cruciale pour améliorer la sécurité, la rentabilité, les performances et la durabilité. AWS fournit un large éventail de mesures liées aux performances à l'aide de services de surveillance tels qu'[Amazon CloudWatch](#) pour vous fournir des informations précieuses. Des indicateurs tels que CPU l'utilisation, l'utilisation de la mémoire, les E/S du disque et les entrées et sorties du réseau peuvent fournir des informations sur les niveaux d'utilisation ou les goulots d'étranglement des performances. Utilisez ces métriques dans le cadre d'une approche fondée sur les données pour ajuster activement et optimiser les ressources de votre charge de travail. Dans un scénario idéal, vous devriez collecter toutes les métriques relatives à vos ressources de calcul sur une plateforme unique, avec des stratégies de conservation mises en œuvre pour atteindre les objectifs financiers et opérationnels.

Étapes d'implémentation

- Identifiez les métriques liées aux performances qui sont pertinentes pour votre charge de travail. Vous devriez collecter des métriques relatives à l'utilisation des ressources et au fonctionnement de votre charge de travail cloud (comme le temps de réponse et le débit).
 - [Métriques EC2 par défaut d'Amazon](#)
 - [Métriques ECS par défaut d'Amazon](#)
 - [Métriques EKS par défaut d'Amazon](#)
 - [Métriques par défaut de Lambda](#)
 - [Métriques relatives à EC2 la mémoire et au disque Amazon](#)
- Choisissez et configurez la solution de journalisation et de surveillance adaptée à votre charge de travail.
 - [Observabilité native AWS](#)
 - [AWS Distro pour OpenTelemetry](#)
 - [Amazon Managed Service for Prometheus](#)
- Définissez le filtre et l'agrégation requis pour les métriques en fonction de vos exigences en matière de charge de travail.
 - [Quantifiez les métriques personnalisées des applications avec Amazon CloudWatch Logs et les filtres métriques](#)
 - [Collectez des statistiques personnalisées grâce au balisage CloudWatch stratégique d'Amazon](#)

- Configurez des stratégies de conservation des données pour vos métriques afin qu'elles correspondent à vos objectifs sécuritaires et opérationnels.
 - [Conservation des données par défaut pour les CloudWatch métriques](#)
 - [Conservation des données par défaut pour les CloudWatch journaux](#)
- Si nécessaire, créez des alarmes et des notifications pour vos métriques afin de vous aider à résoudre de manière proactive les problèmes liés aux performances.
 - [Créez des alarmes pour des métriques personnalisées à l'aide de la détection des CloudWatch anomalies Amazon](#)
 - [Créez des métriques et des alarmes pour des pages Web spécifiques avec Amazon CloudWatch RUM](#)
- Utilisez l'automatisation pour déployer vos agents d'agrégation de métriques et de journaux.
 - [AWS Systems Manager automatisation](#)
 - [OpenTelemetryCollectionneur](#)

Ressources

Documents connexes :

- [Surveillance et observabilité](#)
- [Bonnes pratiques : mise en œuvre de l'observabilité avec AWS](#)
- [CloudWatch Documentation Amazon](#)
- [Collectez des métriques et des journaux à partir d'EC2instances Amazon et de serveurs sur site avec l'agent CloudWatch](#)
- [Accès à Amazon CloudWatch Logs pour AWS Lambda](#)
- [Utilisation CloudWatch des journaux avec des instances de conteneur](#)
- [Publier des métriques personnalisées](#)
- [AWS Réponse : journalisation centralisée](#)
- [AWS Services qui publient CloudWatch des métriques](#)
- [Surveillance d'Amazon EKS sur AWS Fargate](#)

Vidéos connexes :

- [AWS re:Invent 2023 — \[LAUNCH\] Surveillance des applications pour les charges de travail modernes](#)
- [AWS re:Invent 2023 — Mise en œuvre de l'observabilité des applications](#)
- [AWS re:Invent 2023 — Élaborer une stratégie d'observabilité efficace](#)
- [AWS re:Invent 2023 — Une observabilité sans faille avec Distro pour AWS OpenTelemetry](#)
- [Gestion des performances des applications sur AWS](#)

Exemples connexes :

- [AWS Journée d'immersion pour les charges de travail Linux - Amazon CloudWatch](#)
- [Surveillance des ECS clusters et des conteneurs Amazon](#)
- [Surveillance à l'aide des tableaux de CloudWatch bord Amazon](#)
- [EKSAtelier Amazon](#)

PERF02-BP04 Configurer et dimensionner correctement les ressources de calcul

Configurez et dimensionnez correctement les ressources de calcul en fonction des exigences de performance de votre charge de travail et évitez de sous-utiliser ou de surexploiter les ressources.

Anti-modèles courants :

- Vous ignorez les exigences de performance de votre charge de travail, ce qui entraîne un surprovisionnement ou un sous-provisionnement des ressources de calcul.
- Vous ne choisissez que la plus grande ou la plus petite instance disponible pour toutes les charges de travail.
- Vous n'utilisez qu'une seule famille d'instances pour faciliter la gestion.
- Vous ignorez les recommandations de AWS Cost Explorer Compute Optimizer concernant le dimensionnement correct.
- Vous ne réévaluez pas la charge de travail pour voir si de nouveaux types d'instances pourraient convenir.
- Vous ne certifiez qu'un petit nombre de configurations d'instance pour votre organisation.

Avantages liés au respect de cette bonne pratique : dimensionner correctement les ressources de calcul garantit le fonctionnement optimal dans le cloud en évitant le surprovisionnement et le sous-provisionnement des ressources. Le dimensionnement correct des ressources de calcul se traduit généralement par des performances renforcées, une meilleure expérience client et une baisse des coûts.

Niveau de risque encouru si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Le dimensionnement correct permet aux organisations d'exploiter leur infrastructure cloud de manière efficace et rentable tout en répondant aux besoins de l'entreprise. Le surprovisionnement des ressources cloud peut entraîner des coûts supplémentaires, tandis que le sous-provisionnement peut entraîner des performances médiocres et une expérience client négative. AWS fournit des outils tels que [AWS Compute Optimizer](#) et [AWS Trusted Advisor](#) qui utilisent des données historiques pour fournir des recommandations afin de dimensionner correctement vos ressources informatiques.

Étapes d'implémentation

- Choisissez le type d'instance qui correspond le mieux à vos besoins :
 - [Comment choisir le type d'EC2instance Amazon adapté à ma charge de travail ?](#)
 - [Sélection du type d'instance basée sur les attributs pour Amazon Fleet EC2](#)
 - [Créer un groupe Auto Scaling en utilisant la sélection du type d'instance basée sur des attributs](#)
 - [Optimisation de vos coûts de calcul Kubernetes avec la consolidation Karpenter](#)
- Analysez les différentes caractéristiques de performance de votre charge de travail et le lien entre ces caractéristiques et la mémoire, le réseau et CPU l'utilisation. Utilisez ces données pour choisir les ressources qui correspondent le mieux aux objectifs de votre charge de travail en matière de profil et de performance.
- Surveillez l'utilisation de vos ressources à l'aide d'outils de AWS surveillance tels qu'Amazon CloudWatch.
- Sélectionnez la configuration adaptée à vos ressources de calcul.
 - Pour les charges de travail éphémères, évaluez les [CloudWatch indicateurs Amazon](#) de l'instance, CPUUtilization afin de déterminer si l'instance est sous-utilisée ou surutilisée.
 - Pour des charges de travail stables, vérifiez les AWS outils de redimensionnement tels que AWS Compute Optimizer et AWS Trusted Advisor à intervalles réguliers pour identifier les opportunités d'optimisation et de dimensionnement correct de la ressource de calcul.

- Testez les changements de configuration dans un environnement hors production avant de les implémenter dans un environnement réel.
- Réévaluez en permanence les nouvelles offres de calcul et comparez-les aux besoins de votre charge de travail.

Ressources

Documents connexes :

- [Cloud Compute avec AWS](#)
- [Types d'EC2instances Amazon](#)
- [Amazon ECS Containers : instances de ECS conteneurs Amazon](#)
- [EKSConteneurs Amazon : Amazon EKS Worker Nodes](#)
- [Fonctions : configuration des fonctions Lambda](#)
- [Contrôle de l'état du processeur pour votre EC2 instance Amazon](#)

Vidéos connexes :

- [EC2Fondations Amazon](#)
- [AWS re:Invent 2023 — AWS Graviton : le meilleur rapport qualité/prix pour vos charges de travail AWS](#)
- [AWS re:Invent 2023 — Nouvelles fonctionnalités d'IA EC2 générative d'Amazon dans AWS Management Console](#)
- [AWS re:Invent 2023 — Nouveautés d'Amazon EC2](#)
- [AWS re:Invent 2023 — Économies intelligentes : stratégies d'optimisation des coûts d'Amazon EC2](#)
- [AWS re:Invent 2021 — Au service d'EC2Amazon de nouvelle génération : étude approfondie du système Nitro](#)
- [AWS re:Invent 2019 — Amazon Foundations EC2](#)

Exemples connexes :

- [AWS Compute Optimizer Code de démonstration](#)
- [EKSAtelier Amazon](#)
- [Recommandations en matière de redimensionnement](#)

PERF02-BP05 Adaptez dynamiquement vos ressources informatiques

Utilisez l'élasticité du cloud pour mettre à l'échelle vos ressources de calcul de manière dynamique afin de répondre à vos besoins et d'éviter de surprovisionner ou de sous-provisionner la capacité de votre charge de travail.

Anti-modèles courants :

- Vous réagissez aux alertes en augmentant manuellement la capacité.
- Vous utilisez les mêmes recommandations de dimensionnement (généralement, infrastructure statique) que sur site.
- Vous conservez une capacité accrue après un événement de mise à l'échelle au lieu de la réduire.

Avantages liés au respect de cette bonne pratique : en configurant et en testant l'élasticité des ressources de calcul, vous pouvez économiser de l'argent, maintenir les points de référence des performances et améliorer la fiabilité en fonction de l'évolution du trafic.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

AWS offre la flexibilité nécessaire pour augmenter ou diminuer vos ressources de manière dynamique grâce à divers mécanismes de mise à l'échelle afin de répondre à l'évolution de la demande. Combinée aux métriques liées au calcul, la mise à l'échelle dynamique permet aux charges de travail de réagir automatiquement aux changements et d'utiliser l'ensemble optimal de ressources de calcul pour atteindre son objectif.

Vous pouvez utiliser plusieurs approches pour rapprocher l'offre de ressources de la demande.

- Approche de suivi des objectifs : surveillez votre métrique de capacité de mise à l'échelle et augmentez ou réduisez automatiquement votre capacité selon vos besoins.
- Mise à l'échelle prédictive : mettez à l'échelle en prévision des tendances quotidiennes et hebdomadaires.
- Approche basée sur le calendrier : définissez votre propre calendrier de mise à l'échelle en fonction de changements de charge prévisibles.

- Mise à l'échelle des services : choisissez des services (sans serveur, par exemple) conçus pour se mettre à l'échelle automatiquement.

Vous devez vous assurer que les déploiements de charge de travail peuvent gérer les événements de mise à l'échelle ascendante et descendante.

Étapes d'implémentation

- Les instances de calcul, les conteneurs et les fonctions fournissent des mécanismes d'élasticité, soit en combinaison avec l'autoscaling, soit en tant que fonctionnalité du service. Voici des exemples de mécanismes d'autoscaling :

Mécanisme d'autoscaling	Où utiliser
Amazon EC2 Auto Scaling	Pour vous assurer que vous disposez du nombre correct d'EC2 instances Amazon disponibles pour gérer la charge utilisateur de votre application.
Application Autoscaling	Pour dimensionner automatiquement les ressources pour des AWS services individuelles autres qu'Amazon, EC2 tels que AWS Lambda les fonctions ou les services Amazon Elastic Container Service (Amazon ECS) .
Outil Cluster Autoscaler/Karpenter de Kubernetes	Pour mettre à l'échelle automatiquement les clusters Kubernetes.

- La mise à l'échelle est souvent évoquée en lien avec les services de calcul tels que EC2 les instances ou AWS Lambda les fonctions Amazon. Assurez-vous également de prendre en compte la configuration des services non liés au calcul tels que [AWS Glue](#) pour répondre à la demande.
- Vérifiez que les métriques de mise à l'échelle correspondent aux caractéristiques de la charge de travail en cours de déploiement. Si vous déployez une application de transcodage vidéo, un taux d'CPU utilisation de 100 % est attendu et ne doit pas être votre indicateur principal. Utilisez plutôt la profondeur de la file d'attente des tâches de transcodage. Le cas échéant, vous pouvez utiliser une [métrique personnalisée](#) pour votre politique de dimensionnement. Pour choisir les bons indicateurs, prenez en compte les conseils suivants destinés à Amazon EC2 :

- La métrique doit être une métrique d'utilisation valide et décrire à quel point l'instance est occupée.
- La valeur de métrique doit augmenter ou diminuer en proportion du nombre d'instances présentes dans le groupe Auto Scaling.
- Assurez-vous d'utiliser une mise à [l'échelle dynamique](#) plutôt qu'une [mise à l'échelle manuelle](#) pour votre groupe Auto Scaling. Nous vous recommandons également d'utiliser des [politiques de dimensionnement pour le suivi des cibles](#) dans votre dimensionnement dynamique.
- Vérifiez que les déploiements de charges de travail peuvent gérer les deux événements de mise à l'échelle (augmentation et diminution des charges de travail). Par exemple, vous pouvez utiliser [l'historique des activités pour vérifier une activité](#) de mise à l'échelle dans un groupe Auto Scaling.
- Évaluez votre charge de travail pour les modèles prédictifs et mettez-la à l'échelle de manière proactive pour anticiper les changements prévisibles et prévus de la demande. Avec la mise à l'échelle prédictive, vous pouvez supprimer le besoin de surprovisionner de la capacité. Pour plus de détails, consultez [Predictive Scaling with Amazon EC2 Auto Scaling](#).

Ressources

Documents connexes :

- [Cloud Compute avec AWS](#)
- [Types d'EC2instances Amazon](#)
- [Amazon ECS Containers : instances de ECS conteneurs Amazon](#)
- [EKSConteneurs Amazon : Amazon EKS Worker Nodes](#)
- [Fonctions : configuration des fonctions Lambda](#)
- [Contrôle de l'état du processeur pour votre EC2 instance Amazon](#)
- [Présentation approfondie d'Amazon ECS Cluster Auto Scaling](#)
- [Présentation de Karpenter, un Cluster Autoscaler de Kubernetes hautement performant et open source](#)

Vidéos connexes :

- [AWS re:Invent 2023 — AWS Graviton : le meilleur rapport qualité/prix pour vos charges de travail AWS](#)

- [AWS re:Invent 2023 — Nouvelles fonctionnalités d'IA EC2 générative d'Amazon dans Management Console AWS](#)
- [AWS re:Invent 2023 — Nouveautés d'Amazon EC2](#)
- [AWS re:Invent 2023 — Économies intelligentes : stratégies d'optimisation des coûts d'Amazon EC2](#)
- [AWS re:Invent 2021 — Au service d'EC2 Amazon de nouvelle génération : étude approfondie du système Nitro](#)
- [AWS re:Invent 2019 — Amazon Foundations EC2](#)

Exemples connexes :

- [Exemples de groupes Amazon EC2 Auto Scaling](#)
- [EKSAtelier Amazon](#)
- [Faites évoluer vos EKS charges de travail Amazon en exécutant sur IPv6](#)

PERF02-BP06 Utiliser des accélérateurs de calcul matériels optimisés

Utilisez des accélérateurs matériels pour exécuter certaines fonctions de manière plus efficace que les alternatives basées sur l'UC.

Anti-modèles courants :

- En ce qui concerne votre charge de travail, vous n'avez pas comparé une instance à usage général à une instance dédiée qui est capable de fournir de meilleures performances à moindre coût.
- Vous utilisez des accélérateurs de calcul matériels pour les tâches qui peuvent être plus efficaces en utilisant des alternatives basées sur l'UC.
- Vous ne surveillez pas l'utilisation du GPU.

Avantages liés au respect de cette bonne pratique : en utilisant des accélérateurs matériels, tels que des unités de traitement graphique (GPU) et une matrice de portes programmables sur site (FPGA), vous pouvez exécuter certaines fonctions de traitement de manière plus efficace.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Les instances de calcul accéléré donnent accès à des accélérateurs de calcul matériels tels que les GPU et les FPGA. Ces accélérateurs matériels exécutent certaines fonctions comme le traitement graphique ou la correspondance de modèles de données plus efficacement que les alternatives basées sur l'UC. De nombreuses charges de travail accélérées, telles que le rendu, le transcodage et le machine learning, sont très variables en matière d'utilisation des ressources. Exécutez ce matériel uniquement pendant le temps nécessaire et mettez-le hors service grâce à l'automatisation lorsque vous n'en avez plus besoin afin d'améliorer l'efficacité globale des performances.

Étapes d'implémentation

- Identifiez les [instances de calcul accéléré](#) qui peuvent répondre à vos besoins.
- Pour les charges de travail de machine learning, tirez parti d'un matériel conçu spécialement pour votre charge de travail, par exemple [AWS, Trainium](#), [AWS Inferentia](#) et [Amazon EC2 DL1](#). AWS Les instances Inferentia telles que les instances Inf2 [offrent des performances/watt jusqu'à 50 % supérieures à celles des instances Amazon EC2 comparables](#).
- Collectez des métriques d'utilisation pour vos instances de calcul accéléré. Par exemple, vous pouvez utiliser l'agent CloudWatch pour collecter des métriques telles que `utilization_gpu` et `utilization_memory` pour vos GPU, comme indiqué dans [Collecter les métriques des GPU NVIDIA avec Amazon CloudWatch](#).
- Optimisez le code, le fonctionnement du réseau et les paramètres des accélérateurs matériels pour veiller à ce que le matériel sous-jacent soit pleinement utilisé.
 - [Optimiser les paramètres GPU](#)
 - [Surveillance et optimisation des GPU dans l'AMI Deep Learning](#)
 - [Optimisation des E/S pour le réglage des performances de GPU pour l'entraînement du deep learning dans l'IA Amazon SageMaker](#)
- Utilisez les dernières bibliothèques performantes et les pilotes GPU.
- Utilisez l'automatisation pour libérer les instances GPU lorsqu'elles ne sont pas utilisées.

Ressources

Documents connexes :

- [Fonctionnement d'Amazon Elastic Container Service](#)

- [Instances GPU](#)
- [Instances avec AWS Trainium](#)
- [Instances avec AWS Inferentia](#)
- [Passons à l'architecture ! Architecture avec des puces personnalisées et des accélérateurs](#)

- [Calcul accéléré](#)
- [Instances Amazon EC2 VT1](#)
- [Comment choisir le type d'instance EC2 approprié pour ma charge de travail ?](#)
- [Choix du meilleur accélérateur d'IA et de la meilleure compilation de modèles pour l'inférence de vision par ordinateur avec l'IA Amazon SageMaker](#)

Vidéos connexes :

- [AWSre:Invent 2021 - Comment sélectionner les instances Amazon Elastic Compute Cloud GPU pour le deep learning](#)
- [AWSre:INVENT 2022 - \[NOUVEAU LANCEMENT !\] Présentation des instances AWS Amazon EC2 Inf2 basées sur Inferentia2](#)
- [AWSre:Invent 2022 - Accélérez le deep learning et innovez plus rapidement avec Trainium AWS](#)
- [AWS re:Invent 2022 - Deep learning on AWS with NVIDIA: From training to deployment](#)

Exemples connexes :

- [IA Amazon SageMaker et NVIDIA GPU Cloud \(NGC\)](#)
- [Utilisation de l'IA SageMaker avec Trainium et Inferentia pour optimiser les charges de travail d'inférence et d'entraînement du deep learning](#)
- [Optimisation des modèles NLP avec les instances Amazon Elastic Compute Cloud Inf1 dans l'IA Amazon SageMaker](#)

Gestion des données

La solution optimale de gestion des données pour un système particulier varie en fonction du type de données (bloc, fichier ou objet), des modèles d'accès (aléatoire ou séquentiel), du débit requis, de la fréquence d'accès (en ligne, hors-ligne, archivage), de la fréquence de mise à jour (WORM, dynamique), ainsi que des contraintes de disponibilité et de durabilité. Les charges de travail Well-Architected utilisent des magasins de données sur mesure qui intègrent différentes fonctionnalités pour améliorer les performances.

Ce domaine d'intérêt partage des conseils et de bonnes pratiques pour optimiser le stockage de données, les modèles de déplacement et d'accès, ainsi que l'efficacité des performances des magasins de données.

Bonnes pratiques

- [PERF03-BP01 Utilisez un magasin de données spécialement conçu pour répondre au mieux à vos besoins en matière d'accès aux données et de stockage](#)
- [PERF03-BP02 Évaluer les options de configuration disponibles pour le magasin de données](#)
- [PERF03-BP03 Collecter et enregistrer les indicateurs de performance du magasin de données](#)
- [PERF03-BP04 Mise en œuvre de stratégies pour améliorer les performances des requêtes dans un magasin de données](#)
- [PERF03-BP05 Implémenter des modèles d'accès aux données qui utilisent la mise en cache](#)

PERF03-BP01 Utilisez un magasin de données spécialement conçu pour répondre au mieux à vos besoins en matière d'accès aux données et de stockage

Comprenez les caractéristiques des données (telles que la possibilité de partage, la taille, la taille du cache, les modèles d'accès, la latence, le débit et la persistance des données) afin de sélectionner les magasins de données dédiés (stockage ou base de données) adaptés à votre charge de travail.

Anti-modèles courants :

- Vous vous en tenez à un magasin de données, car l'équipe interne sait comment tirer parti de ce type de solution en particulier.

- Vous partez du principe que toutes les charges de travail ont des exigences similaires en termes de stockage de données et d'accès aux données.
- Vous n'avez pas implémenté de catalogue de données pour inventorier vos ressources de données.

Avantages liés au respect de cette bonne pratique : en comprenant l'importance des caractéristiques et des exigences des données, vous pouvez déterminer la technologie de stockage la plus efficace et la plus performante adaptée à vos besoins en matière de charge de travail.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

Lors de la sélection et de la mise en œuvre du stockage des données, assurez-vous que les caractéristiques de requête, de dimensionnement et de stockage répondent aux exigences relatives aux données de charge de travail. AWS fournit de nombreuses technologies de stockage de données et de base de données, notamment le stockage par blocs, le stockage d'objets, le stockage en continu, les systèmes de fichiers, les bases de données relationnelles, les bases de données à valeur clé, les bases de données documentaires, en mémoire, les graphiques, les séries chronologiques et les bases de données de registre. Chaque solution de gestion de données propose des options et des configurations pour prendre en charge vos cas d'utilisation et vos modèles de données. En comprenant les caractéristiques et les exigences des données, vous pouvez vous affranchir de la technologie de stockage monolithique et des one-size-fits-all approches restrictives pour vous concentrer sur la gestion appropriée des données.

Étapes d'implémentation

- Procédez à l'inventaire des différents types de données qui existent dans votre charge de travail.
- Comprenez et documentez les caractéristiques et les exigences des données, notamment :
 - Type de données (non structurées, semi-structurées, relationnelles)
 - Volume et croissance des données
 - Durabilité des données : persistantes, éphémères, temporaires
 - ACIDexigences (atomicité, consistance, isolation, durabilité)
 - Modèles d'accès aux données (à lecture intensive ou à écriture intensive)
 - Latence
 - Débit

- IOPS(opérations d'entrée/sortie par seconde)
- Période de conservation des données
- Découvrez les différents magasins de données (services [de stockage](#) et [de base](#) de données) disponibles pour votre charge de travail AWS qui peuvent répondre aux caractéristiques de vos données, comme indiqué dans [PERF01-BP01 Découvrez et comprenez les services et fonctionnalités cloud disponibles](#). Voici quelques exemples de technologies de stockage AWS et leurs principales caractéristiques :

Type	AWS Services	Principales caractéristiques
Stockage d'objets	Amazon S3	Capacité de mise à l'échelle illimitée, haute disponibilité et plusieurs options d'accessibilité. Le transfert et l'accès à des objets à l'intérieur et à l'extérieur d'Amazon S3 peuvent utiliser un service, tel que Transfer Acceleration ou Access Points (points d'accès), pour répondre à votre localisation, à vos besoins en matière de sécurité et à vos modèles d'accès.
Archivage et stockage	Amazon S3 Glacier	Conçu pour l'archivage des données.
Stockage en streaming	Amazon Kinesis Amazon Managed Streaming pour Apache Kafka (Amazon MSK)	Ingestion et stockage efficaces des données de streaming.
Système de fichiers partagé	Amazon Elastic File System (AmazonEFS)	Système de fichiers montable auquel plusieurs types de

Type	AWS Services	Principales caractéristiques
		solutions informatiques peuvent accéder.
Système de fichiers partagé	Amazon FSx	Construit sur les dernières solutions AWS informatiques pour prendre en charge quatre systèmes de fichiers couramment utilisés : Open NetApp ONTAPZFS, Windows File Server et Lustre. FSx La latence, le débit et le débit d'Amazon IOPS varient selon le système de fichiers et doivent être pris en compte lors de la sélection du système de fichiers adapté à vos besoins en matière de charge de travail.
Stockage en mode bloc	Boutique Amazon Elastic Block (AmazonEBS)	Service de stockage par blocs évolutif et performant conçu pour Amazon Elastic Compute Cloud (AmazonEC2). Amazon EBS inclut le stockage SSD sauvegardé pour les charges de travail transactionnelles intensives et le stockage HDD sauvegardé pour les charges de travail IOPS gourmandes en débit.

Type	AWS Services	Principales caractéristiques
Base de données relationnelle	Amazon Aurora , Amazon RDS , Amazon Redshift .	Conçu pour prendre en charge les transactions ACID (atomicité, cohérence, isolation, durabilité) et pour maintenir l'intégrité référentielle et la forte cohérence des données. De nombreuses applications traditionnelles, de planification des ressources d'entreprise (ERP), de gestion de la relation client (CRM) et de commerce électronique utilisent des bases de données relationnelles pour stocker leurs données.
Base de données clé-valeur	Amazon DynamoDB	Optimisées pour les modèles d'accès courants, généralement pour stocker et récupérer de gros volumes de données. Les applications Web à trafic élevé, les systèmes d'e-commerce et les applications de jeu sont des cas d'utilisation typiques pour les bases de données de valeurs-clés.

Type	AWS Services	Principales caractéristiques
Base de données documentaire	Amazon DocumentDB	Conçu pour stocker des données semi-structurées sous JSON forme de documents similaires. Ces bases de données aident les développeurs à créer et mettre à jour rapidement des applications telles que la gestion de contenu, les catalogues et les profils utilisateur.
Base de données en mémoire	Amazon ElastiCache , Amazon MemoryDB pour Redis	Utilisées pour les applications qui nécessitent un accès en temps réel aux données, la latence la plus faible et le débit le plus élevé. Vous pouvez utiliser des bases de données en mémoire pour la mise en cache des applications, la gestion des sessions, les classements des jeux, le magasin de fonctionnalités ML à faible latence, le système de messagerie à microservices et un mécanisme de streaming à haut débit

Type	AWS Services	Principales caractéristiques
Base de données orientée graphe	Amazon Neptune	Destinées aux applications qui doivent parcourir et interroger des millions de relations entre des jeux de données graphiques hautement connectés avec une latence de millisecondes à grande échelle. De nombreuses entreprises utilisent des bases de données de graphiques pour la détection des fraudes, les réseaux sociaux et les moteurs de recommandation.
Base de données de séries temporelles	Amazon Timestream	Utilisées pour collecter, synthétiser et extraire efficacement des informations à partir de données qui changent au fil du temps. Les applications IoT et la DevOps télémétrie industrielle peuvent utiliser des bases de données de séries chronologiques.

Type	AWS Services	Principales caractéristiques
Larges colonnes	Amazon Keyspaces (pour Apache Cassandra)	Utilise des tables, des lignes et des colonnes, mais contrairement à une base de données relationnelle, les noms et le format des colonnes peuvent varier d'une ligne à l'autre dans la même table. Généralement, vous voyez un magasin de colonnes larges dans les applications industrielles à grande échelle pour la maintenance des équipements, la gestion des parcs et l'optimisation des itinéraires.
Registre	Base de données Amazon Quantum Ledger (AmazonQLDB)	Fournit une autorité centralisée et fiable pour conserver un enregistrement évolutif, immuable et vérifiable grâce au chiffrement des transactions pour chaque application. Il n'est pas rare de voir des bases de données de registre utilisées pour les systèmes d'enregistrement, la chaîne d'approvisionnement, les inscriptions et même les transactions bancaires.

- Si vous créez une plate-forme de données, tirez parti de [l'architecture de données moderne](#) AWS pour intégrer votre lac de données, votre entrepôt de données et vos magasins de données spécialement conçus.
- Les principales questions que vous devez vous poser lors du choix d'un magasin de données pour votre charge de travail sont les suivantes :

Question	Éléments à prendre en compte
Comment sont structurées les données ?	<ul style="list-style-type: none">• Si les données ne sont pas structurées, envisagez un magasin d'objets tel qu'Amazon S3 ou une base de SQL données sans base de données telle qu'Amazon DocumentDB• Pour les données clé-valeur, pensez à DynamoDB, Amazon (ElasticCache Redis) ou Amazon MemoryDB OSS
Quel niveau d'intégrité référentielle est requis ?	<ul style="list-style-type: none">• En ce qui concerne les contraintes liées aux clés étrangères, les bases de données relationnelles telles qu'Amazon RDS et Aurora peuvent fournir ce niveau d'intégrité.• Généralement, dans un SQL modèle sans données, vous dénormaliseriez vos données en un seul document ou en un ensemble de documents à récupérer en une seule demande plutôt que de joindre plusieurs documents ou tableaux.
La conformité ACID (atomicité, consistance, isolation, durabilité) est-elle requise ?	<ul style="list-style-type: none">• Si les ACID propriétés associées aux bases de données relationnelles sont requises, envisagez une base de données relationnelle telle qu'Amazon et RDS Aurora.• Si une forte cohérence est requise pour Aucune SQL base de données, vous pouvez utiliser des lectures fortement cohérentes avec DynamoDB.

Question	Éléments à prendre en compte
<p>Comment les exigences de stockage vont-elles évoluer au fil du temps ? Comment cela affectera-t-il la capacité de mise à l'échelle ?</p>	<ul style="list-style-type: none"> • Les bases de données sans serveur telles que DynamoDB et Amazon Quantum Ledger Database (QLDBAmazon) évolueront de manière dynamique. • Les bases de données relationnelles ont des limites supérieures sur le stockage alloué et doivent souvent être partitionnées horizontalement à l'aide de mécanismes tels que le partitionnement une fois qu'elles atteignent ces limites.
<p>Quelle est la proportion de requêtes en lecture par rapport aux requêtes en écriture ? La mise en cache pourrait-elle améliorer les performances ?</p>	<ul style="list-style-type: none"> • Les charges de travail gourmandes en lecture peuvent bénéficier d'une couche de mise en cache, comme ElastiCache ou DAX si la base de données est DynamoDB. • Les lectures peuvent également être déchargées pour lire des répliques avec des bases de données relationnelles telles qu'Amazon. RDS
<p>Le stockage et la modification (OLTP- Traitement des transactions en ligne) ou la récupération et le reporting (OLAP- Traitement analytique en ligne) ont-ils une priorité plus élevée ?</p>	<ul style="list-style-type: none"> • Pour un traitement transactionnel en lecture telle quelle à haut débit, envisagez une base de données sans base de données SQL telle que DynamoDB. • Pour des modèles de lecture complexes et à haut débit (comme la jointure) cohérents, utilisez Amazon. RDS • Pour les requêtes analytiques, envisagez d'utiliser une base de données en colonnes telle qu'Amazon Redshift ou d'exporter les données vers Amazon S3 et d'effectuer des analyses à l'aide d'Athena ou d'Amazon. QuickSight

Question	Éléments à prendre en compte
Quel est le niveau de durabilité requis pour les données ?	<ul style="list-style-type: none">• Aurora réplique automatiquement vos données sur trois zones de disponibilité au sein d'une région. Autrement dit, vos données sont très durables avec moins de risque de perte de données.• DynamoDB est automatiquement répliqué sur plusieurs zones de disponibilité, assurant ainsi la haute disponibilité et la durabilité des données.• Amazon S3 offre une durabilité de 99,999999999 %. De nombreux services de base de données, tels qu'Amazon RDS et DynamoDB, prennent en charge l'exportation de données vers Amazon S3 pour une conservation et un archivage à long terme.
Souhaitez-vous vous éloigner des moteurs de base de données commerciaux ou des coûts de licence ?	<ul style="list-style-type: none">• Pensez aux moteurs open source tels que Postgre SQL et My on SQL Amazon ou RDS Aurora.• Tirez parti de AWS Database Migration Service et AWS Schema Conversion Tool pour passer des moteurs de bases de données commerciaux vers des moteurs open source.
Qu'attendez-vous de la base de données du point de vue opérationnel ? Le passage aux services gérés est-il une préoccupation majeure ?	<ul style="list-style-type: none">• Tirer parti d'Amazon RDS au lieu d'AmazonEC2, et de DynamoDB ou d'Amazon DocumentDB au lieu d'héberger vous-même une SQL base de données « No » peut réduire les frais d'exploitation.

Question	Éléments à prendre en compte
<p>Comment accédez-vous actuellement à la base de données ? S'agit-il uniquement d'un accès aux applications ou existe-t-il des utilisateurs de Business Intelligence (BI) et d'autres off-the-shelf applications connectées ?</p>	<ul style="list-style-type: none"> • Si vous dépendez d'outils externes, vous devrez peut-être maintenir la compatibilité avec les bases de données qu'ils prennent en charge. Amazon RDS est entièrement compatible avec les différentes versions de moteurs qu'il prend en charge, notamment Microsoft SQL Server, OracleSQL, My et PostgreSQL.

- Réalisez des tests et procédez au benchmarking dans un environnement hors production afin d'identifier le magasin de données qui peut répondre à vos exigences en termes de charge de travail.

Ressources

Documents connexes :

- [Types de EBS volumes Amazon](#)
- [EC2Stockage Amazon](#)
- [Amazon EFS : Amazon EFS Performance](#)
- [Amazon FSx pour Lustre Performance](#)
- [Performances du serveur de fichiers Amazon FSx pour Windows](#)
- [Amazon Glacier S3 : documentation Amazon Glacier S3](#)
- [Amazon S3 : directives en matière de débit de demandes et de performances](#)
- [Stockage dans le cloud avec AWS](#)
- [Caractéristiques d'Amazon EBS I/O](#)
- [Bases de données cloud avec AWS](#)
- [AWS Mise en cache de bases de données](#)
- [DynamoDB Accelerator](#)
- [Bonnes pratiques Amazon Aurora](#)
- [Performances d'Amazon Redshift](#)
- [Amazon Athena top 10 de conseils en matière de performance](#)

- [Bonnes pratiques Amazon Redshift Spectrum](#)
- [Bonnes pratiques Amazon DynamoDB](#)
- [Choisissez entre Amazon EC2 et Amazon RDS](#)
- [Bonnes pratiques pour la mise en œuvre d'Amazon ElastiCache](#)

Vidéos connexes :

- [AWS re:Invent 2023 : Améliorez l'efficacité d'Amazon Elastic Block Store et soyez plus rentable](#)
- [AWS re:Invent 2023 : Optimisation du prix et des performances du stockage avec Amazon Simple Storage Service](#)
- [AWS re:Invent 2023 : Création et optimisation d'un lac de données sur Amazon Simple Storage Service](#)
- [AWS re:Invent 2022 : Création d'architectures de données modernes sur AWS](#)
- [AWS re:Invent 2022 : Création d'architectures de maillage de données sur AWS](#)
- [AWS re:Invent 2023 : présentation approfondie d'Amazon Aurora et de ses innovations](#)
- [AWS re:Invent 2023 : Modélisation avancée des données avec Amazon DynamoDB](#)
- [AWS re:Invent 2022 : Modernisez les applications avec des bases de données spécialement conçues](#)
- [Amazon DynamoDB deep dive: Advanced design patterns](#)

Exemples connexes :

- [AWS Atelier sur les bases de données spécialement conçues](#)
- [Bases de données pour développeurs](#)
- [AWS Journée d'immersion dans l'architecture de données moderne](#)
- [Créez un maillage de données sur AWS](#)
- [Exemples Amazon S3](#)
- [Optimisation du modèle de données à l'aide du partage de données Amazon Redshift](#)
- [Migrations des bases de données](#)
- [MS SQL Server - AWS Database Migration Service \(AWS DMS\) Démo de réplication](#)
- [Atelier pratique sur la modernisation des bases de données](#)

- [Échantillons Amazon Neptune](#)

PERF03-BP02 Évaluer les options de configuration disponibles pour le magasin de données

Comprenez et évaluez les différentes fonctionnalités et options de configuration disponibles pour vos magasins de données afin d'optimiser l'espace de stockage et les performances de votre charge de travail.

Anti-modèles courants :

- Vous n'utilisez qu'un seul type de stockage, tel qu'AmazonEBS, pour toutes les charges de travail.
- Vous utilisez le provisionné IOPS pour toutes les charges de travail sans effectuer de tests réels sur tous les niveaux de stockage.
- Vous ne connaissez pas les options de configuration de la solution de gestion de données que vous avez choisie.
- Vous vous concentrez uniquement sur l'augmentation de la taille de l'instance sans examiner les autres options de configuration disponibles.
- Vous ne testez pas les caractéristiques de mise à l'échelle de votre magasin de données.

Avantages liés au respect de cette bonne pratique : en explorant et en expérimentant les configurations de magasin de données, vous pourriez réduire le coût de l'infrastructure, améliorer les performances et réduire l'effort requis pour maintenir vos charges de travail.

Niveau de risque encouru si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Une charge de travail peut comporter un ou plusieurs magasins de données utilisés en fonction des exigences de stockage des données et d'accès aux données. Pour optimiser l'efficacité et le coût de vos performances, vous devez évaluer les modèles d'accès aux données afin de déterminer les configurations de magasin de données appropriées. Pendant que vous explorez les options de magasin de données, tenez compte de divers aspects tels que les options de stockage, la mémoire, le calcul, le réplica en lecture, les exigences de cohérence, le regroupement de connexions et les options de mise en cache. Testez ces différentes options de configuration pour améliorer les métriques d'efficacité des performances.

Étapes d'implémentation

- Comprenez les configurations actuelles (comme le type d'instance, la taille de stockage ou la version du moteur de base de données) de votre magasin de données.
- Consultez AWS la documentation et les meilleures pratiques pour découvrir les options de configuration recommandées qui peuvent contribuer à améliorer les performances de votre magasin de données. Les principales options de magasin de données à prendre en compte sont les suivantes :

Option de configuration	Exemples
Déchargement des lectures (comme les réplicas en lecture et la mise en cache)	<ul style="list-style-type: none"> • Pour les tables DynamoDB, vous pouvez décharger les lectures à l'aide de la mise en cache. DAX • Vous pouvez créer un cluster Amazon ElastiCache (RedisOSS) et configurer votre application pour qu'elle lise d'abord dans le cache, puis revenir à la base de données si l'élément demandé n'est pas présent. • Bases de données relationnelles telles qu'Amazon RDS et Aurora, et mises en service Aucune SQL base de données telle que Neptune et Amazon DocumentDB ne prend toutes en charge l'ajout de répliques de lecture pour décharger les parties de lecture de la charge de travail. • Les bases de données sans serveur comme DynamoDB se mettent à l'échelle automatiquement. Assurez-vous de disposer de suffisamment d'unités de capacité de lecture (RCU) pour gérer la charge de travail.
Mise à l'échelle des écritures (comme le partitionnement des clés de partition ou l'introduction d'une file d'attente)	<ul style="list-style-type: none"> • Pour les bases de données relationnelles, vous pouvez augmenter la taille de l'instance pour faire face à une charge de travail accrue ou augmenter le provisionnement

Option de configuration	Exemples
	<p>IOPs pour augmenter le débit du stockage sous-jacent.</p> <ul style="list-style-type: none">• Vous pouvez également ajouter une file d'attente devant votre base de données plutôt que d'écrire directement dans la base de données. Ce modèle vous permet de dissocier l'ingestion de la base de données et de contrôler le débit afin que la base de données ne soit pas submergée.• Regrouper vos demandes d'écriture plutôt que de créer de nombreuses transactions de courte durée contribue à améliorer le débit dans les bases de données relationnelles à volume d'écriture élevé.• Les bases de données sans serveur telles que DynamoDB peuvent augmenter le débit d'écriture automatiquement ou en ajustant les unités de capacité d'écriture allouées WCU () en fonction du mode de capacité.• Vous pouvez toujours rencontrer des problèmes avec les partitions à chaud lorsque vous atteignez les limites de débit pour une clé de partition donnée. Pour pallier ce problème, choisissez une clé de partition distribuée plus uniformément ou partitionnez en écriture la clé de partition.

Option de configuration	Exemples
Politiques de gestion du cycle de vie de vos jeux de données	<ul style="list-style-type: none"> Vous pouvez utiliser Amazon S3 Lifecycle afin de gérer vos objets au cours de leur cycle de vie. Si vos schémas d'accès sont inconnus, changeants ou imprévisibles, vous pouvez utiliser Amazon S3 Intelligent-Tiering, qui surveille les schémas d'accès et déplace automatiquement les objets qui n'ont pas été accédés vers des niveaux d'accès moins coûteux. Vous pouvez tirer parti des métriques Amazon S3 Storage Lens pour identifier les opportunités d'optimisation et les lacunes dans la gestion du cycle de vie. Amazon EFS Lifecycle Management gère automatiquement le stockage de fichiers pour vos systèmes de fichiers.
Gestion et regroupement des connexions	<ul style="list-style-type: none"> Amazon RDS Proxy peut être utilisé avec Amazon RDS et Aurora pour gérer les connexions à la base de données. Les bases de données sans serveur comme DynamoDB n'ont pas de connexions associées, mais tenez compte de la capacité allouée et des stratégies de mise à l'échelle automatique pour faire face aux pics de charge.

- Réalisez des tests et procédez au benchmarking dans un environnement hors production afin d'identifier l'option de configuration qui répond à vos exigences en termes de charge de travail.
- Après avoir réalisé vos tests, planifiez votre migration et validez vos métriques de performance.
- Utilisez AWS des outils de surveillance (comme [Amazon CloudWatch](#)) et d'optimisation (comme [Amazon S3 Storage Lens](#)) pour optimiser en permanence votre magasin de données en utilisant des modèles d'utilisation réels.

Ressources

Documents connexes :

- [Stockage cloud avec AWS](#)
- [Types de EBS volumes Amazon](#)
- [EC2Stockage Amazon](#)
- [Amazon EFS : Amazon EFS Performance](#)
- [Amazon FSx pour Lustre Performance](#)
- [Performances du serveur de fichiers Amazon FSx pour Windows](#)
- [Amazon Glacier S3 : documentation Amazon Glacier S3](#)
- [Amazon S3 : directives en matière de débit de demandes et de performances](#)
- [Caractéristiques d'Amazon EBS I/O](#)
- [Bases de données cloud avec AWS](#)
- [AWS Mise en cache de bases de données](#)
- [DynamoDB Accelerator](#)
- [Bonnes pratiques Amazon Aurora](#)
- [Performances d'Amazon Redshift](#)
- [Amazon Athena top 10 de conseils en matière de performance](#)
- [Bonnes pratiques Amazon Redshift Spectrum](#)
- [Bonnes pratiques Amazon DynamoDB](#)

Vidéos connexes :

- [AWS re:Invent 2023: Improve Amazon Elastic Block Store efficiency and be more cost-efficient](#)
- [AWS re:Invent 2023: Optimize storage price and performance with Amazon Simple Storage Service](#)
- [AWS re:Invent 2023: Building and optimizing a data lake on Amazon Simple Storage Service](#)
- [AWS re:Invent 2023 : Nouveautés en matière de stockage de fichiers AWS](#)
- [AWS re:Invent 2023: Dive deep into Amazon DynamoDB](#)

Exemples connexes :

- [AWS Atelier sur les bases de données spécialement conçues](#)
- [Bases de données pour développeurs](#)
- [AWS Journée d'immersion dans l'architecture de données moderne](#)
- [Amazon EBS Autoscale](#)
- [Exemples Amazon S3](#)
- [Exemple Amazon DynamoDB](#)
- [AWS Exemples de migration de base de données](#)
- [Atelier sur la modernisation des bases de données](#)
- [Utilisation des paramètres de votre base de données Amazon RDS pour Postgress](#)

PERF03-BP03 Collecter et enregistrer les indicateurs de performance du magasin de données

Suivez et archivez les métriques de performance pertinentes pour votre magasin de données afin de comprendre comment fonctionnent vos solutions de gestion des données. Ces métriques peuvent vous aider à optimiser votre magasin de données, à vérifier que les exigences de votre charge de travail sont satisfaites et à fournir une vue d'ensemble claire sur le fonctionnement de la charge de travail.

Anti-modèles courants :

- Vous utilisez uniquement la recherche manuelle des fichiers journaux pour les métriques.
- Vous publiez uniquement des métriques sur les outils internes utilisés par votre équipe et vous n'avez pas une visibilité complète de votre charge de travail.
- Vous n'utilisez que les métriques par défaut enregistrées par le logiciel de surveillance que vous avez sélectionné.
- Vous n'examinez les métriques qu'en cas de problème.
- Vous ne surveillez que les métriques au niveau du système et vous ne capturez pas les métriques d'accès aux données ou d'utilisation des données.

Avantages liés au respect de cette bonne pratique : la définition de points de référence pour les performances vous permet de mieux comprendre le comportement normal et les exigences des charges de travail. Les modèles anormaux peuvent être identifiés et débogués plus rapidement, ce qui améliore les performances et la fiabilité du magasin de données.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

L'enregistrement de plusieurs métriques de performance sur une période donnée est nécessaire pour la surveillance des performances de vos magasins de données. Cette surveillance vous permet non seulement de détecter les anomalies, mais aussi d'évaluer les performances par rapport aux métriques métier afin de vérifier que vous répondez aux besoins de votre charge de travail.

Ces métriques doivent inclure à la fois le système sous-jacent qui prend en charge le magasin de données et les métriques de la base de données. Les indicateurs système sous-jacents peuvent inclure CPU l'utilisation, la mémoire, le stockage sur disque disponible, les E/S sur disque, le taux d'accès au cache et les mesures entrantes et sortantes du réseau, tandis que les indicateurs du magasin de données peuvent inclure les transactions par seconde, les requêtes les plus fréquentes, les taux de requêtes moyens, les temps de réponse, l'utilisation de l'index, les blocages de table, les délais d'attente des requêtes et le nombre de connexions ouvertes. Ces données sont essentielles pour comprendre comment fonctionne la charge de travail et comment la solution de gestion des données est utilisée. Utilisez ces métriques dans le cadre d'une approche fondée sur les données pour ajuster et optimiser les ressources de votre charge de travail.

Utilisez des outils, des bibliothèques et des systèmes qui enregistrent des mesures de performances liées aux performances de la base de données.

Étapes d'implémentation

- Identifiez les métriques de performances clés que votre magasin de données doit suivre.
 - [Métriques et dimensions d'Amazon S3](#)
 - [Mesures de surveillance pour une RDS instance Amazon](#)
 - [Surveillance de la charge de base de données avec Performance Insights sur Amazon RDS](#)
 - [Vue d'ensemble de la surveillance améliorée](#)
 - [Métriques et dimensions DynamoDB](#)
 - [Surveillance de l'accélérateur DynamoDB](#)
 - [Surveillance d'Amazon MemoryDB avec Amazon CloudWatch](#)
 - [Quelles métriques dois-je surveiller ?](#)
 - [Surveillance des performances de cluster Amazon Redshift](#)
 - [Métriques et dimensions Timestream](#)

- [CloudWatch Métriques Amazon pour Amazon Aurora](#)
- [Journalisation et surveillance dans Amazon Keyspaces \(pour Apache Cassandra\)](#)
- [Surveillance des ressources Amazon Neptune](#)
- Utilisez une solution de journalisation et de surveillance approuvée pour collecter ces métriques. [Amazon CloudWatch](#) peut collecter des métriques sur l'ensemble des ressources de votre architecture. Vous pouvez également récupérer et publier des métriques personnalisées pour faire apparaître des métriques d'entreprise ou des métriques dérivées. Utilisez CloudWatch ou utilisez des solutions tierces pour définir des alarmes indiquant lorsque les seuils sont dépassés.
- Vérifiez si la surveillance du magasin de données peut bénéficier d'une solution de machine learning qui détecte les anomalies de performance.
 - [Amazon DevOps Guru for Amazon RDS](#) fournit de la visibilité sur les problèmes de performance et recommande des mesures correctives.
- Configurez la conservation des données dans votre solution de surveillance et de journalisation en fonction de vos objectifs sécuritaires et opérationnels.
 - [Conservation des données par défaut pour les CloudWatch métriques](#)
 - [Conservation des données par défaut pour les CloudWatch journaux](#)

Ressources

Documents connexes :

- [Mise en cache de bases de données AWS](#)
- [Amazon Athena top 10 de conseils en matière de performance](#)
- [Bonnes pratiques Amazon Aurora](#)
- [DynamoDB Accelerator](#)
- [Bonnes pratiques Amazon DynamoDB](#)
- [Bonnes pratiques Amazon Redshift Spectrum](#)
- [Performances d'Amazon Redshift](#)
- [Bases de données cloud avec AWS](#)
- [Amazon RDS Performance Insights](#)

Vidéos connexes :

- [AWS re:Invent 2022 - Surveillance des performances avec Amazon et RDS Aurora, avec Autodesk](#)
- [Surveillance et optimisation des performances des bases de données avec Amazon DevOps Guru pour Amazon RDS](#)
- [AWS re:Invent 2023 - Nouveautés en matière de stockage de fichiers AWS](#)
- [AWS re:Invent 2023 - Découvrez en détail Amazon DynamoDB](#)
- [AWS re:Invent 2023 - Création et optimisation d'un lac de données sur Amazon S3](#)
- [AWS re:Invent 2023 - Nouveautés en matière de stockage de fichiers AWS](#)
- [AWS re:Invent 2023 - Découvrez en détail Amazon DynamoDB](#)
- [Meilleures pratiques pour surveiller les charges de travail Redis sur Amazon ElastiCache](#)

Exemples connexes :

- [Cadre de collecte de métriques pour l'ingestion des jeux de données AWS](#)
- [Atelier RDS de surveillance Amazon](#)
- [AWS Atelier sur les bases de données spécialement conçues](#)

PERF03-BP04 Mise en œuvre de stratégies pour améliorer les performances des requêtes dans un magasin de données

Mettez en œuvre des stratégies pour optimiser les données et améliorer les requêtes sur les données afin de renforcer la capacité de mise à l'échelle et l'efficacité des performances pour votre charge de travail.

Anti-modèles courants :

- Vous ne partitionnez pas les données dans votre magasin de données.
- Vous ne stockez les données que dans un seul format de fichier dans votre magasin de données.
- Vous n'utilisez pas d'index dans votre magasin de données.

Avantages liés au respect de cette bonne pratique : en optimisant les performances des données et des requêtes, vous augmentez leur efficacité, vous réduisez les coûts et vous améliorez l'expérience utilisateur.

Niveau de risque encouru si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

L'optimisation des données et des requêtes sont des aspects essentiels de l'efficacité des performances d'un magasin de données, car ils ont un impact sur les performances et la réactivité de l'ensemble de la charge de travail dans le cloud. Les données non optimisées peuvent augmenter l'utilisation des ressources et les goulots d'étranglement, ce qui réduit l'efficacité globale d'un magasin de données.

L'optimisation des données inclut plusieurs techniques pour garantir un stockage de données et un accès aux données efficaces. Cela permet également d'améliorer les performances des requêtes dans un magasin de données. Les principales stratégies incluent le partitionnement des données, la compression des données et la dénormalisation des données, qui permettent d'optimiser les données à la fois pour le stockage et l'accès.

Étapes d'implémentation

- Comprenez et analysez les requêtes essentielles sur les données effectuées dans votre magasin de données.
- Identifiez les requêtes lentes dans votre magasin de données et utilisez des plans de requêtes pour comprendre leur état actuel.
 - [Analyse du plan de requêtes dans Amazon Redshift](#)
 - [Utilisation d'EXPLAIN et EXPLAIN ANALYZE sur Athena](#)
- Mettez en œuvre des stratégies pour améliorer les performances des requêtes. Les stratégies clés incluent :
 - L'utilisation d'un [format de fichier en colonnes](#) (comme Parquet ou ORC).
 - La compression des données dans le magasin de données pour réduire l'espace de stockage et les opérations d'E/S.
 - Le partitionnement des données pour diviser les données en parties plus petites et réduire le temps d'analyse des données.
 - [Partitionnement de données dans Athena](#)
 - [Partitions et distribution des données](#)
 - L'indexation des données sur les colonnes communes de la requête.
 - Utilisez des vues matérialisées pour les requêtes fréquentes.
 - [Compréhension des vues matérialisées](#)
 - [Création de vues matérialisées dans Amazon Redshift](#)

- Choisissez l'opération de jointure appropriée pour la requête. Lorsque vous joignez deux tables, spécifiez la table la plus grande sur le côté gauche de la jointure et la plus petite sur le côté droit de la jointure.
- La solution de mise en cache distribué pour améliorer la latence et réduire le nombre d'opérations d'E/S dans la base de données.
- Maintenance régulière, telle que l'[aspiration](#), la réindexation et les [statistiques d'exécution](#).
- Expérimentez et testez les stratégies dans un environnement hors production.

Ressources

Documents connexes :

- [Bonnes pratiques Amazon Aurora](#)
- [Performances d'Amazon Redshift](#)
- [Amazon Athena top 10 de conseils en matière de performance](#)
- [Mise en cache de bases de données AWS](#)
- [Bonnes pratiques de mise en œuvre d'Amazon ElastiCache](#)
- [Partitionnement de données dans Athena](#)

Vidéos connexes :

- [AWS re:Invent 2023 - AWS storage cost-optimization best practices](#)
- [AWS re:Invent 2022 - Performance monitoring with Amazon RDS and Aurora, featuring Autodesk](#)
- [Optimize Amazon Athena Queries with New Query Analysis Tools](#)

Exemples connexes :

- [AWS Atelier sur les bases de données sur mesure](#)

PERF03-BP05 Implémenter des modèles d'accès aux données qui utilisent la mise en cache

Mettez en œuvre des modèles d'accès qui peuvent tirer parti de la mise en cache des données pour une récupération rapide des données fréquemment consultées.

Anti-modèles courants :

- Vous mettez en cache des données qui changent fréquemment.
- Vous utilisez les données mises en cache comme si elles étaient stockées de manière durable et toujours disponibles.
- Vous ne tenez pas compte de la cohérence de vos données mises en cache.
- Vous ne surveillez pas l'efficacité de la mise en cache.

Avantages liés au respect de cette bonne pratique : le stockage des données dans un cache contribue à améliorer la latence et le débit de lecture, l'expérience utilisateur et l'efficacité globale, tout en réduisant les coûts.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Un cache est un composant logiciel ou matériel destiné à stocker des données afin que les requêtes futures portant sur les mêmes données puissent être traitées plus rapidement ou plus efficacement. Les données stockées dans un cache peuvent être reconstruites en cas de perte en répétant un calcul antérieur ou en les récupérant dans un autre magasin de données.

La mise en cache des données peut être l'une des stratégies les plus efficaces pour améliorer les performances globales de votre application et réduire la charge qui pèse sur vos sources de données principales sous-jacentes. Les données peuvent être mises en cache à plusieurs niveaux de l'application, par exemple au sein de l'application en effectuant des appels à distance ou mise en cache côté client ou en utilisant un service secondaire rapide pour stocker les données mise en cache à distance.

Mise en cache côté client

Grâce à la mise en cache côté client, chaque client (une application ou un service qui interroge l'entrepôt de données dorsales) peut stocker les résultats de ses requêtes uniques localement pendant une durée spécifiée. Cela permet de réduire le nombre de requêtes adressées à un entrepôt de données sur le réseau en vérifiant d'abord le cache du client local. En l'absence de résultats, l'application peut alors interroger l'entrepôt de données et stocker ces résultats localement. Ce modèle permet à chaque client de stocker les données dans l'emplacement le plus proche possible (le client lui-même), ce qui se traduit par la latence la plus faible possible. Les clients peuvent

également continuer à répondre à certaines requêtes lorsque l'entrepôt de données dorsales n'est pas disponible, ce qui augmente la disponibilité de l'ensemble du système.

L'un des inconvénients de cette approche est que lorsque plusieurs clients sont impliqués, ils peuvent stocker les mêmes données mises en cache localement. Cela entraîne à la fois une double utilisation du stockage et une incohérence des données entre ces clients. Un client peut mettre en cache les résultats d'une requête et, une minute plus tard, un autre client peut exécuter la même requête et obtenir un résultat différent.

Mise en cache à distance

Pour résoudre le problème de duplication de données entre clients, un service externe rapide ou un cache distant, peut être utilisé pour stocker les données demandées. Au lieu de vérifier un magasin de données local, chaque client vérifie le cache distant avant d'interroger l'entrepôt de données dorsales. Cette stratégie permet d'obtenir des réponses plus cohérentes entre les clients, d'améliorer l'efficacité des données stockées et d'augmenter le volume de données mises en cache, car l'espace de stockage évolue indépendamment des clients.

L'inconvénient d'un cache distant est que l'ensemble du système peut connaître une latence plus élevée, car un saut de réseau à réseau supplémentaire est nécessaire pour vérifier le cache distant. La mise en cache côté client peut être utilisée parallèlement à la mise en cache à distance pour une mise en cache à plusieurs niveaux afin d'améliorer la latence.

Étapes d'implémentation

- Identifiez les bases de données APIs et les services réseau susceptibles de bénéficier de la mise en cache. Les services dont la charge de travail de lecture est importante, dont le read-to-write ratio est élevé ou dont la mise à l'échelle est coûteuse sont candidats à la mise en cache.
 - [Mise en cache de bases de données](#)
 - [Activation de la API mise en cache pour améliorer la réactivité](#)
- Identifiez le type de stratégie de mise en cache le mieux adapté à votre modèle d'accès.
 - [Stratégies de mise en cache](#)
 - [Solutions de mise en cache AWS](#)
- Suivez les [bonnes pratiques de mise en cache](#) pour votre banque de données.
- Configurez une stratégie d'invalidation du cache, telle que a time-to-live (TTL), pour toutes les données afin d'équilibrer la fraîcheur des données et de réduire la pression sur la banque de données principale.

- Activez des fonctionnalités telles que les nouvelles tentatives de connexion automatiques, le backoff exponentiel, les délais d'attente côté client et le regroupement des connexions dans le client, le cas échéant, car elles peuvent améliorer les performances et la fiabilité.
 - [Bonnes pratiques : clients Redis et Amazon ElastiCache \(RedisOSS\)](#)
- Surveillez le taux d'accès au cache en visant un objectif de 80 % ou plus. Des valeurs inférieures peuvent indiquer une taille de cache insuffisante ou un modèle d'accès qui ne bénéficie pas de la mise en cache.
 - [Quelles métriques dois-je surveiller ?](#)
 - [Bonnes pratiques pour surveiller les charges de travail Redis sur Amazon ElastiCache](#)
 - [Surveillance des meilleures pratiques avec Amazon ElastiCache \(RedisOSS\) à l'aide d'Amazon CloudWatch](#)
- Mettre en œuvre la [réplication des données](#) pour transférer les lectures vers plusieurs instances et améliorer les performances et la disponibilité de lecture des données.

Ressources

Documents connexes :

- [Utilisation de l'objectif Amazon ElastiCache Well-Architected](#)
- [Surveillance des meilleures pratiques avec Amazon ElastiCache \(RedisOSS\) à l'aide d'Amazon CloudWatch](#)
- [Quelles métriques dois-je surveiller ?](#)
- [ElastiCache Livre blanc sur les performances à grande échelle avec Amazon](#)
- [Défis et stratégies en matière de mise en cache](#)

Vidéos connexes :

- [Parcours de ElastiCache formation Amazon](#)
- [Concevez pour réussir grâce aux ElastiCache meilleures pratiques d'Amazon](#)
- [AWS re:Invent 2020 - Concevez pour réussir grâce aux meilleures pratiques d'Amazon ElastiCache](#)
- [AWS re:Invent 2023 - \[LAUNCH\] Présentation d'Amazon Serverless ElastiCache](#)
- [AWS re:Invent 2022 - 5 excellentes façons de réinventer votre couche de données avec Redis](#)
- [AWS re:Invent 2021 - Présentation approfondie d'Amazon ElastiCache \(Redis\) OSS](#)

Exemples connexes :

- [Améliorer les performances SQL de ma base de données avec Amazon ElastiCache \(RedisOSS\)](#)

Réseau et diffusion de contenu

La solution de mise en réseau optimale pour une charge de travail varie en fonction de la latence, des exigences de débit, de l'instabilité et de la bande passante. Le choix des options d'emplacement est tributaire des contraintes physiques telles que les ressources pour utilisateur ou sur site. Ces contraintes peuvent être compensées avec les emplacements périphériques ou le placement des ressources.

Sur AWS, la mise en réseau est virtualisée et disponible dans plusieurs types et configurations. Il est ainsi plus facile de répondre à vos besoins en matière de réseau. AWS propose des fonctionnalités de produit (par exemple, Enhanced Networking, Amazon EC2 networking optimized instances, Amazon S3 transfer acceleration et Amazon CloudFront dynamique) pour optimiser le trafic réseau. AWS propose également des fonctionnalités de mise en réseau (par exemple, Amazon Route 53 latency routing, des points de terminaison Amazon VPC, AWS Direct Connect et AWS Global Accelerator) pour réduire la distance ou la gigue du réseau.

Ce domaine d'intérêt partage des conseils et de bonnes pratiques pour concevoir, configurer et exploiter des solutions de mise en réseau et de diffusion de contenu efficaces dans le cloud.

Bonnes pratiques

- [PERF04-BP01 Comprendre l'impact du réseau sur les performances](#)
- [PERF04-BP02 Évaluer les fonctionnalités réseau disponibles](#)
- [PERF04-BP03 Choisissez une connectivité dédiée adaptée à votre charge VPN de travail](#)
- [PERF04-BP04 Utiliser l'équilibrage de charge pour répartir le trafic entre plusieurs ressources](#)
- [PERF04-BP05 Choisissez les protocoles réseau pour améliorer les performances](#)
- [PERF04-BP06 Choisissez l'emplacement de votre charge de travail en fonction des exigences du réseau](#)
- [PERF04-BP07 Optimiser la configuration du réseau en fonction des métriques](#)

PERF04-BP01 Comprendre l'impact du réseau sur les performances

Analysez et comprenez l'impact des décisions liées au réseau sur votre charge de travail afin de fournir des performances efficaces et une meilleure expérience utilisateur.

Anti-modèles courants :

- Tout le trafic passe par vos centres de données existants.
- Vous acheminez l'ensemble du trafic via des pare-feux centralisés au lieu d'utiliser des outils de sécurité réseau natifs cloud.
- Vous configurez AWS Direct Connect des connexions sans connaître les exigences d'utilisation réelles.
- Vous ne tenez pas compte des caractéristiques de la charge de travail et de la surcharge de chiffrement lors de la définition de vos solutions de mise en réseau.
- Vous utilisez des concepts et des stratégies sur site pour les solutions de mise en réseau dans le cloud.

Avantages liés au respect de cette bonne pratique : comprendre comment la mise en réseau affecte les performances de la charge de travail vous aide à identifier les goulots d'étranglement potentiels, à améliorer l'expérience utilisateur, à accroître la fiabilité et à réduire la maintenance opérationnelle à mesure que la charge de travail évolue.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

Le réseau est responsable de la connectivité entre les composants d'application, les services cloud, les réseaux périphériques et les données sur site et, par conséquent, il peut avoir un impact majeur sur les performances de la charge de travail. Outre les performances de la charge de travail, l'expérience utilisateur peut également être affectée par la latence du réseau, la bande passante, les protocoles, l'emplacement, la congestion du réseau, l'instabilité, le débit et les règles de routage.

Veillez à avoir une liste documentée des exigences de mise en réseau de la charge de travail, y compris la latence, la taille des paquets, les règles de routage, les protocoles et les modèles de trafic pris en charge. Passez en revue les solutions de mise en réseau disponibles et identifiez le service qui répond aux caractéristiques de mise en réseau de votre charge de travail. Les réseaux basés sur le cloud peuvent être rapidement recréés. L'évolution de votre architecture réseau au fil du temps est donc nécessaire pour améliorer l'efficacité des performances.

Étapes d'implémentation :

- Définissez et documentez les exigences de performance réseau, y compris les métriques telles que la latence du réseau, la bande passante, les protocoles, les emplacements, les modèles de trafic (pics et fréquence), le débit, le chiffrement, l'inspection et les règles de routage.
- Découvrez les principaux services AWS réseau tels que [VPCs](#), [Elastic Load Balancing \(ELB\)](#) et [Amazon Route 53](#). [AWS Direct Connect](#)
- Capturez les principales caractéristiques réseau suivantes :

Caractéristiques	Outils et métriques
Caractéristiques de mise en réseau fondamentales	<ul style="list-style-type: none"> • VPC Journaux de flux • AWS Transit Gateway Journaux de flux • AWS Transit Gateway métriques • AWS PrivateLink métriques
Caractéristiques de mise en réseau des applications	<ul style="list-style-type: none"> • Elastic Fabric Adapter (EFA) • AWS App Mesh métriques • Métriques Amazon API Gateway
Caractéristiques de mise en réseau à la périphérie	<ul style="list-style-type: none"> • CloudFront Métriques Amazon • Métriques Amazon Route 53 • AWS Global Accelerator métriques
Caractéristiques de mise en réseau hybride	<ul style="list-style-type: none"> • AWS Direct Connect métriques • AWS Site-to-Site VPN métriques • AWS Client VPN métriques • AWS Cloud WAN métriques
Caractéristiques de mise en réseau de la sécurité	<ul style="list-style-type: none"> • AWS Shield, AWS WAF, et AWS Network Firewall métriques
Caractéristiques de traçage	<ul style="list-style-type: none"> • AWS X-Ray • VPC Analyseur de Reachability • Analyseur d'accès réseau

Caractéristiques	Outils et métriques
	<ul style="list-style-type: none">• Amazon Inspector• Amazon CloudWatch RUM

- Définition de points de référence et test des performances du réseau :
 - [Comparez](#) le débit du réseau, car certains facteurs peuvent affecter les performances EC2 du réseau Amazon lorsque les instances se trouvent dans les mêmes VPC instances. Mesurez la bande passante réseau entre les instances Amazon EC2 Linux d'une même instanceVPC.
 - Effectuez des [tests de charge](#) pour expérimenter des solutions et des options de mise en réseau.

Ressources

Documents connexes :

- [Application Load Balancer](#)
- [EC2Mise en réseau améliorée sous Linux](#)
- [EC2Mise en réseau améliorée sous Windows](#)
- [EC2Groupes de placement](#)
- [Activation de la mise en réseau améliorée avec l'adaptateur réseau Elastic \(ENA\) sur les instances Linux](#)
- [Network Load Balancer](#)
- [Produits de mise en réseau avec AWS](#)
- [Passerelle de transit](#)
- [Transition vers le routage basé sur la latence dans Amazon Route 53](#)
- [VPCPoints de terminaison](#)

Vidéos connexes :

- [AWS re:Invent 2023 - AWS mise en réseau des fondations](#)
- [AWS re:Invent 2023 - Que peut apporter le réseau à votre application ?](#)
- [AWS re:Invent 2023 - VPC Designs avancés et nouvelles fonctionnalités](#)
- [AWS re:Invent 2023 - Guide du développeur sur les réseaux cloud](#)
- [AWS re:Invent 2019 - Connectivité AWS et architectures réseau hybrides AWS](#)

- [AWS re:Invent 2019 - Optimisation des performances réseau pour les instances Amazon EC2](#)
- [AWS Summit Online - Améliorez les performances du réseau mondial pour les applications](#)
- [AWS re:Invent 2020 - Meilleures pratiques et astuces de mise en réseau avec le Well-Architected Framework](#)
- [AWS re:Invent 2020 : meilleures pratiques de AWS mise en réseau pour les migrations à grande échelle](#)

Exemples connexes :

- [AWS Transit Gateway et solutions de sécurité évolutives](#)
- [AWS Ateliers de réseautage](#)
- [Atelier pratique sur le pare-feu réseau](#)
- [Observation et diagnostic de votre réseau sur AWS](#)
- [Recherche et résolution des erreurs de configuration réseau sur AWS](#)

PERF04-BP02 Évaluer les fonctionnalités réseau disponibles

Évaluez les fonctions de mise en réseau dans le cloud qui peuvent améliorer les performances. Mesurez l'impact de ces fonctions au moyen de tests, de métriques et de l'analyse. Par exemple, tirez parti des fonctionnalités au niveau du réseau qui sont disponibles pour réduire la latence, la distance réseau ou l'instabilité.

Anti-modèles courants :

- Vous restez au sein d'une même région, car c'est là que votre siège social se trouve physiquement.
- Vous utilisez des pare-feux plutôt que des groupes de sécurité pour filtrer le trafic.
- Vous faites une pause TLS pour inspecter le trafic plutôt que de vous fier aux groupes de sécurité, aux politiques relatives aux terminaux et à d'autres fonctionnalités natives du cloud.
- Vous utilisez uniquement la segmentation basée sur un sous-réseau au lieu des groupes de sécurité.

Avantages liés au respect de cette bonne pratique : l'évaluation de toutes les options et fonctionnalités de service peut augmenter les performances de vos charges de travail, baisser le coût d'infrastructure, réduire les efforts nécessaires à la maintenance de vos charges de travail et

améliorer votre posture générale en matière de sécurité. Vous pouvez utiliser le AWS backbone mondial pour offrir une expérience réseau optimale à vos clients.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

AWS propose des services tels [AWS Global Accelerator](#) CloudFront qu'[Amazon](#) qui peuvent contribuer à améliorer les performances du réseau, tandis que la plupart AWS des services proposent des fonctionnalités (telles que la fonctionnalité [Amazon S3 Transfer Acceleration](#)) permettant d'optimiser le trafic réseau.

Examinez les options de configuration liées au réseau disponibles et leur impact potentiel sur votre charge de travail. L'optimisation des performances dépend de la compréhension de la manière dont ces options interagissent avec votre architecture et de l'impact qu'elles auront à la fois sur les performances mesurées et sur l'expérience utilisateur.

Étapes d'implémentation

- Créer une liste des composants de la charge de travail.
 - Envisagez [AWS Cloud WAN](#) de l'utiliser pour créer, gérer et surveiller le réseau de votre organisation lors de la création d'un réseau mondial unifié.
 - Surveillez vos réseaux mondiaux et principaux à l'aide des [métriques Amazon CloudWatch Logs](#). Tirez parti d'[Amazon CloudWatch RUM](#), qui fournit des informations permettant d'identifier, de comprendre et d'améliorer l'expérience numérique des utilisateurs.
 - Visualisez la latence réseau globale entre les zones de disponibilité Régions AWS et au sein de chaque zone de disponibilité, [AWS Network Manager](#) afin de mieux comprendre le lien entre les performances de votre application et les performances du AWS réseau sous-jacent.
 - Utilisez un outil ou un service de base de données de gestion de configuration (CMDB) existant, par exemple [AWS Config](#) pour créer un inventaire de votre charge de travail et de sa configuration.
- Identifier et documenter le test comparatif pour vos métriques de performances s'il s'agit d'une charge de travail existante, en vous concentrant sur les goulots d'étranglement et les zones à améliorer. Les métriques de mise en réseau liées aux performances diffèrent par charge de travail en fonction des exigences métier et des caractéristiques de charge de travail. Pour commencer, il pourrait être important d'examiner ces métriques pour votre charge de travail : bande passante, latence, perte de paquets, instabilité et retransmissions.

- S'il s'agit d'une nouvelle charge de travail, effectuez des [tests de charge](#) pour identifier les goulots d'étranglement liés aux performances.
- Concernant l'identification des goulots d'étranglement au niveau des performances, examiner les options de configuration pour les solutions afin d'identifier les opportunités d'amélioration des performances. Découvrez les principales options et fonctionnalités de mise en réseau suivantes :

Opportunité d'amélioration	Solution
Chemin ou itinéraires réseau	Utilisez l' analyseur d'accès réseau pour identifier les chemins ou les itinéraires.
Protocoles réseau	Consultez PERF04-BP05 Choisissez les protocoles réseau pour améliorer les performances .
Topologie du réseau	<p>Évaluez vos compromis opérationnels et de performance entre le VPCpeering et AWS Transit Gateway lors de la connexion de plusieurs comptes. AWS Transit Gateway simplifie la façon dont vous interconnectez tous vos VPCs réseaux, qui peuvent s'étendre sur Comptes AWS des milliers de réseaux locaux. Partagez votre compte AWS Transit Gateway entre plusieurs comptes en utilisant AWS Resource Access Manager.</p> <p>Consultez PERF04-BP03 Choisissez une connectivité dédiée adaptée à votre charge VPN de travail.</p>
Services de réseau	<p>AWS Global Accelerator est un service réseau qui améliore les performances du trafic de vos utilisateurs jusqu'à 60 % en utilisant l'infrastructure réseau AWS mondiale.</p> <p>Amazon CloudFront peut améliorer les performances de votre charge de travail, de</p>

Opportunité d'amélioration	Solution
	<p>diffusion de contenu et de latence à l'échelle mondiale.</p> <p>Utilisez Lambda @edge pour exécuter des fonctions qui personnalisent le contenu au plus près CloudFront des utilisateurs, réduisent le temps de latence et améliorent les performances.</p> <p>Amazon Route 53 propose des options de routage basées sur la latence, de routage de géolocalisation, de routage de géolocalisation et de routage basé sur IP pour vous aider à améliorer les performances de votre charge de travail auprès d'un public mondial. Identifiez l'option de routage qui optimiserait les performances de votre charge de travail en examinant le trafic de votre charge de travail et la localisation des utilisateurs lorsque votre charge de travail est distribuée dans le monde entier.</p>

Opportunité d'amélioration	Solution
Fonctionnalités des ressources de stockage	<p>Amazon S3 Transfer Acceleration est une fonctionnalité qui permet aux utilisateurs externes de bénéficier des optimisations du réseau CloudFront pour télécharger des données vers Amazon S3. Cela améliore le transfert d'importants volumes de données à partir d'emplacements distants qui n'ont pas de connectivité dédiée au AWS Cloud.</p> <p>Les points d'accès multi-régions Amazon S3 répliquent le contenu vers plusieurs régions et simplifient la charge de travail en fournissant un point d'accès. Lorsqu'un point d'accès multi-région est utilisé, vous pouvez demander ou écrire des données à Amazon S3 tandis que le service identifie le compartiment à la latence la plus faible.</p>

Opportunité d'amélioration	Solution
Fonctionnalités des ressources informatiques	<p>Les interfaces réseau élastiques (ENA) utilisées par EC2 les instances Amazon, les conteneurs et les fonctions Lambda sont limitées par flux. Passez en revue vos groupes de placement pour optimiser le débit EC2 de votre réseau. Pour éviter un goulot d'étranglement par flux, créez votre application pour qu'elle utilise plusieurs flux. Pour surveiller et obtenir une visibilité sur vos métriques réseau liées au calcul, utilisez CloudWatch Metrics et ethtool. La <code>ethtool</code> commande est incluse dans le ENA pilote et expose des métriques supplémentaires liées au réseau qui peuvent être publiées sous forme de métrique personnalisée sur CloudWatch</p> <p>Les Amazon Elastic Network Adapters (ENA) fournissent une optimisation supplémentaire en fournissant un meilleur débit à vos instances au sein d'un groupe de placement de clusters.</p> <p>Elastic Fabric Adapter (EFA) est une interface réseau pour les EC2 instances Amazon qui vous permet d'exécuter des charges de travail nécessitant des niveaux élevés de communications entre nœuds à grande échelle. AWS</p> <p>Les instances EBS optimisées pour Amazon utilisent une pile de configuration optimisée et fournissent une capacité dédiée supplémentaire pour augmenter les EBS E/S Amazon.</p>

Ressources

Documents connexes :

- [Application Load Balancer](#)
- [EC2 Mise en réseau améliorée sous Linux](#)
- [EC2 Mise en réseau améliorée sous Windows](#)
- [EC2 Groupes de placement](#)
- [Activation de la mise en réseau améliorée avec l'adaptateur réseau Elastic \(ENA\) sur les instances Linux](#)
- [Network Load Balancer](#)
- [Produits de mise en réseau avec AWS](#)
- [Transition vers le routage basé sur la latence dans Amazon Route 53](#)
- [VPC Points de terminaison](#)
- [Journaux de flux VPC](#)

Vidéos connexes :

- [AWS re:Invent 2023 — Prêts pour la suite ? Designing networks for growth and flexibility](#)
- [AWS re:Invent 2023 — VPC Designs avancés et nouvelles fonctionnalités](#)
- [AWS re:Invent 2023 — Guide du développeur sur les réseaux cloud](#)
- [AWS re:Invent 2022 — Approfondissez l'infrastructure réseau AWS](#)
- [AWS re:Invent 2019 — Connectivité AWS et architectures réseau hybrides AWS](#)
- [AWS re:Invent 2018 — Optimisation des performances réseau pour les instances Amazon EC2](#)
- [AWS Global Accelerator](#)

Exemples connexes :

- [AWS Transit Gateway et solutions de sécurité évolutives](#)
- [AWS Ateliers de réseautage](#)
- [Observation et diagnostic de votre réseau](#)
- [Détection et correction des erreurs de configuration du réseau sur AWS](#)

PERF04-BP03 Choisissez une connectivité dédiée adaptée à votre charge VPN de travail

Lorsque la connectivité hybride est requise pour connecter des ressources sur site et dans le cloud, allouez une bande passante adéquate pour répondre à vos exigences de performance. Estimez les exigences en matière de bande passante et de latence pour votre charge de travail hybride. Ces chiffres détermineront vos exigences en matière de dimensionnement.

Anti-modèles courants :

- Vous évaluez uniquement les VPN solutions en fonction des exigences de chiffrement de votre réseau.
- Vous n'évaluez pas les options de sauvegarde ni de connectivité redondante.
- Vous n'identifiez pas toutes les exigences de la charge de travail (chiffrement, protocole, bande passante et trafic requis).

Avantages liés au respect de cette bonne pratique : la sélection et la configuration de solutions de connectivité appropriées renforcent la fiabilité de votre charge de travail et optimisent les performances. En identifiant les exigences en matière de charge de travail, en planifiant à l'avance et en évaluant les solutions hybrides, vous pouvez minimiser les modifications coûteuses du réseau physique et les frais d'exploitation tout en augmentant votre time-to-value

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

Développez une architecture réseau hybride en fonction de vos besoins en bande passante. [AWS Direct Connect](#) vous permet de connecter votre réseau sur site en privé à AWS. Cette solution convient lorsque vous avez besoin d'une bande passante élevée et d'une faible latence tout en conservant des performances constantes. Une VPN connexion établit une connexion sécurisée via Internet. Elle sert uniquement lorsque seule une connexion temporaire est requise, lorsque le coût est un facteur, ou en cas d'urgence en attendant qu'une connectivité réseau physique résiliente soit établie lors de l'utilisation d' AWS Direct Connect.

Si vos besoins en bande passante sont élevés, vous pouvez envisager plusieurs VPN services AWS Direct Connect ou services. Le trafic peut être équilibré entre les services, mais nous ne recommandons pas l'équilibrage de charge entre AWS Direct Connect et en VPN raison des différences de latence et de bande passante.

Étapes d'implémentation

- Évaluez les besoins en bande passante et en latence de vos applications existantes.
 - Pour les charges de travail existantes qui sont transférées AWS, exploitez les données de vos systèmes de surveillance réseau internes.
 - Pour les nouvelles charges de travail ou pour les charges de travail existantes pour lesquelles vous ne disposez pas de données de suivi, contactez les propriétaires du produit pour obtenir des métriques de performance adéquates et offrir une bonne expérience utilisateur.
- Sélectionnez une connexion dédiée ou VPN comme option de connectivité. En fonction de toutes les exigences en matière de charge de travail (chiffrement, bande passante et besoins en trafic), vous pouvez choisir AWS Direct Connect ou [AWS VPN](#) (ou les deux). Le schéma suivant peut vous aider à choisir le type de connexion approprié.
 - [AWS Direct Connect](#) fournit une connectivité dédiée à l'environnement AWS, de 50 Mbit/s à 100 Gbit/s, en utilisant des connexions dédiées ou des connexions hébergées. Cela vous permet de gérer et de contrôler la latence et de profiter d'une bande passante provisionnée. Ainsi, vos charges de travail peuvent se connecter efficacement à d'autres environnements. En faisant appel à des AWS Direct Connect partenaires, vous pouvez bénéficier d'une end-to-end connectivité à partir de plusieurs environnements, fournissant ainsi un réseau étendu aux performances constantes. AWS permet de dimensionner la bande passante de connexion directe en utilisant 100 Gbit/s natifs, un groupe d'agrégation de liens (LAG) ou un multipath à BGP coût égal (). ECMP
 - AWS [Site-to-Site VPN](#) fournit un VPN service géré prenant en charge la sécurité du protocole Internet (IPsec). Lorsqu'une VPN connexion est créée, chaque VPN connexion inclut deux tunnels pour une haute disponibilité.
- Suivez AWS la documentation pour choisir l'option de connectivité appropriée :
 - Si vous décidez de l'utiliser AWS Direct Connect, sélectionnez la bande passante adaptée à votre connectivité.
 - Si vous utilisez un réseau AWS Site-to-Site VPN sur plusieurs sites pour vous connecter à un Région AWS, utilisez une [Site-to-Site VPN connexion accélérée](#) afin d'améliorer les performances du réseau.
 - Si la conception de votre réseau consiste en IPsec VPN une connexion via une connexion [AWS Direct Connect](#), pensez à utiliser une adresse IP privée VPN pour améliorer la sécurité et réaliser une segmentation. [AWS Site-to-Site L'adresse IP privée VPN](#) est déployée au-dessus de l'interface virtuelle de transit (VIF).

- [AWS Direct Connect SiteLink](#) permet de créer des connexions redondantes et à faible latence entre vos centres de données du monde entier en envoyant les données sur le chemin le plus rapide entre les [AWS Direct Connect sites](#), en les contournant. Régions AWS
- Validez votre configuration de connectivité avant le déploiement en production. Effectuez des tests de sécurité et de performance pour vous assurer qu'elle répond à vos exigences en matière de bande passante, de fiabilité, de latence et de conformité.
- Surveillez régulièrement les performances et l'utilisation de votre connectivité et optimisez-les si nécessaire.

Organigramme des performances déterministes

Ressources

Documents connexes :

- [Produits de mise en réseau avec AWS](#)
- [AWS Transit Gateway](#)
- [VPC Points de terminaison](#)
- [Création d'une infrastructure VPC AWS multiréseau évolutive et sécurisée](#)
- [Client VPN](#)

Vidéos connexes :

- [AWS re:Invent 2023 — Création d'une connectivité réseau hybride avec AWS](#)
- [AWS re:Invent 2023 — Connectivité à distance sécurisée pour AWS](#)
- [AWS re:Invent 2022 — Optimisation des performances avec Amazon CloudFront](#)
- [AWS re:Invent 2019 — Connectivité AWS et architectures réseau hybrides AWS](#)
- [AWS re:Invent 2020 — Connect AWS Transit Gateway](#)

Exemples connexes :

- [AWS Transit Gateway et solutions de sécurité évolutives](#)

- [AWS Ateliers de réseautage](#)

PERF04-BP04 Utiliser l'équilibrage de charge pour répartir le trafic entre plusieurs ressources

Répartissez le trafic sur plusieurs ressources ou services pour permettre à votre charge de travail de tirer parti de l'élasticité fournie par le cloud. Vous pouvez également utiliser l'équilibrage de charge afin de décharger la terminaison du chiffrement en vue d'améliorer les performances, d'assurer la fiabilité et de gérer et acheminer efficacement le trafic.

Anti-modèles courants :

- Vous ne tenez pas compte des exigences de votre charge de travail lorsque vous choisissez le type d'équilibreur de charge.
- Vous ne tirez pas parti des fonctionnalités de l'équilibreur de charge pour optimiser les performances.
- La charge de travail est exposée directement à Internet sans équilibreur de charge.
- Vous acheminez tout le trafic Internet via des équilibreurs de charge existants.
- Vous utilisez un équilibrage de TCP charge générique et vous faites en sorte que chaque nœud de calcul gère SSL le chiffrement.

Avantages liés au respect de cette bonne pratique : un équilibreur de charge gère la charge variable du trafic de votre application dans une seule zone de disponibilité ou entre plusieurs zones de disponibilité et permet une haute disponibilité, une mise à l'échelle automatique et une meilleure utilisation de votre charge de travail.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

Les équilibreurs de charge constituent le point d'entrée de votre charge de travail, à partir duquel ils distribuent le trafic vers vos cibles principales, telles que les instances de calcul ou les conteneurs, afin d'améliorer l'utilisation.

Le choix du bon type d'équilibreur de charge est la première étape de l'optimisation de votre architecture. Commencez par répertorier les caractéristiques de votre charge de travail, telles que

le protocole (comme TCPHTTP, TLS, ou WebSockets), le type de cible (comme les instances, les conteneurs ou les applications sans serveur), les exigences de l'application (telles que les connexions de longue durée, l'authentification des utilisateurs ou la rigidité) et le placement (par exemple, région, zone locale, avant-poste ou isolation zonale).

AWS fournit plusieurs modèles permettant à vos applications d'utiliser l'équilibrage de charge.

[Application Load Balancer](#) convient parfaitement à l'équilibrage de la charge HTTP et du HTTPS trafic et fournit un routage avancé des demandes destiné à la fourniture d'architectures d'applications modernes, notamment des microservices et des conteneurs.

[Network Load Balancer](#) convient parfaitement à l'équilibrage de charge du TCP trafic lorsque des performances extrêmes sont requises. Il est capable de traiter des millions de requêtes par seconde tout en maintenant de très faibles latences. Il est optimisé pour gérer les tendances soudaines et instables du trafic.

[Elastic Load Balancing](#) intègre la gestion et le SSL TLS déchiffrement des certificats, ce qui vous permet de gérer de manière centralisée les SSL paramètres de l'équilibreur de charge et de décharger les tâches CPU intensives de votre charge de travail.

Après avoir choisi le bon équilibreur de charge, vous pouvez commencer à tirer parti de ses fonctionnalités pour réduire les efforts que votre système dorsal doit fournir pour servir le trafic.

Par exemple, en utilisant à la fois Application Load Balancer (ALB) et Network Load Balancer NLB (), vous pouvez SSL effectuer un déchargement par TLS chiffrement, ce qui vous permet d'éviter que vos cibles ne se lancent CPU dans une poignée de main TLS intensive et d'améliorer la gestion des certificats.

Lorsque vous configurez SSL ou TLS déchargez dans votre équilibreur de charge, celui-ci devient responsable du chiffrement du trafic en provenance et à destination des clients, tout en distribuant le trafic non chiffré à vos backends, en libérant les ressources de votre backend et en améliorant le temps de réponse des clients.

Application Load Balancer peut également desservir HTTP /2 trafic sans avoir à le prendre en charge sur vos cibles. Cette simple décision peut améliorer le temps de réponse de votre application, car HTTP /2 utilise TCP les connexions de manière plus efficace.

Les exigences de latence de votre charge de travail doivent être prises en compte lors de la définition de l'architecture. Par exemple, si vous avez une application sensible à la latence, vous pouvez décider d'utiliser Network Load Balancer, qui offre des latences extrêmement faibles. Vous pouvez

également décider de rapprocher votre charge de travail de vos clients en tirant parti d'Application Load Balancer dans [AWS Local Zones](#) or même [AWS Outposts](#).

L'équilibrage de charge entre zones est un autre élément à prendre en compte pour les charges de travail sensibles à la latence. Avec l'équilibrage de charge inter-zone, chaque nœud d'équilibreur de charge distribue le trafic sur les cibles enregistrées dans toutes les zones de disponibilité activées.

Intégrez l'autoscaling à votre équilibreur de charge. L'un des aspects essentiels d'un système performant est le dimensionnement adéquat de vos ressources dorsales. Pour ce faire, vous pouvez tirer parti des intégrations d'équilibreurs de charge pour les ressources cibles du système dorsal. Grâce à l'intégration de l'équilibreur de charge avec les groupes Auto Scaling, les cibles seront ajoutées à l'équilibreur de charge ou retirées de l'équilibreur de charge selon les besoins en fonction du trafic entrant. Les équilibreurs de charge peuvent également s'intégrer à [Amazon ECS et Amazon EKS](#) pour les charges de travail conteneurisées.

- [Amazon ECS - Équilibrage de charge des services](#)
- [Équilibrage de charge des applications sur Amazon EKS](#)
- [Équilibrage de charge réseau sur Amazon EKS](#)

Étapes d'implémentation

- Définissez vos exigences en matière d'équilibrage de charge, notamment en termes de volume de trafic, de disponibilité et de capacité de mise à l'échelle des applications.
- Choisissez le type d'équilibreur de charge adapté à votre application.
 - Utilisez Application Load Balancer pour les charges de travail HTTP/HTTPS.
 - Utilisez Network Load Balancer pour les charges autres que les HTTP charges de travail qui s'exécutent sur ou. TCP UDP
 - Utilisez une combinaison des deux ([ALB comme cible NLB](#)) si vous souhaitez tirer parti des fonctionnalités des deux produits. Par exemple, vous pouvez le faire si vous souhaitez utiliser la statique IPs de NLB avec le routage basé sur les HTTP en-têtes depuis ALB, ou si vous souhaitez exposer votre HTTP charge de travail à un [AWS PrivateLink](#).
- Pour une comparaison complète des équilibreurs de charge, consultez la [comparaison des ELB produits](#).
- Utilisez SSL/TLS offloading si possible.
 - Configurez HTTPS/TLS listeners avec [Application Load Balancer](#) [et Network Load Balancer](#) [intégrés](#) à [AWS Certificate Manager](#)

- Notez que certaines charges de travail peuvent nécessiter un end-to-end chiffrement pour des raisons de conformité. Dans ce cas, il est nécessaire de permettre le chiffrement au niveau des cibles.
- Pour connaître les meilleures pratiques en matière de sécurité, voir [SEC09-BP02 Appliquer le chiffrement](#) en transit.
- Sélectionnez le bon algorithme de routage (uniquement ALB).
 - L'algorithme de routage peut faire une réelle différence dans la manière d'utiliser vos cibles dorsales et donc dans leur impact sur les performances. Par exemple, ALB propose [deux options pour les algorithmes de routage](#) :
 - Demandes en suspens les moins nombreuses : cette option permet d'obtenir une meilleure répartition de la charge sur vos cibles dorsales dans les cas où les requêtes de votre application varient en complexité ou vos cibles varient en capacité de traitement.
 - Tour de rôle : utilisez cette méthode lorsque les requêtes et les cibles sont similaires, ou si vous devez distribuer les requêtes de manière égale entre les cibles.
- Envisagez l'option inter-zone ou l'isolement par zone.
 - Désactivez l'option désactivée (utilisez l'isolement par zone) pour améliorer la latence et les domaines de panne par zone. Il est désactivé par défaut dans NLB et dans, [ALB vous pouvez le désactiver par groupe cible](#).
 - Activez l'option inter-zone pour une disponibilité et une flexibilité accrues. Par défaut, l'interzone est activée pour ALB et [NLB vous pouvez l'activer pour chaque groupe cible](#).
- Activez HTTP Keep-Alives pour vos HTTP charges de travail (uniquement). ALB Grâce à cette fonctionnalité, l'équilibreur de charge peut réutiliser les connexions du backend jusqu'à l'expiration du délai de conservation, améliorant ainsi votre temps de HTTP demande et de réponse et réduisant également l'utilisation des ressources sur vos cibles de backend. Pour plus de détails sur la façon de procéder pour Apache et Nginx, consultez [Quels sont les paramètres optimaux pour utiliser Apache ou en NGINX tant que serveur principal pour ? ELB](#)
- Activez la surveillance pour votre équilibreur de charge.
 - Activez les journaux d'accès pour votre [Application Load Balancer](#) et [Network](#) Load Balancer.
 - Les principaux domaines à prendre en compte ALB sont `request_processing_time`, `request_processing_time`, `response_processing_time`.
 - Les principaux domaines à prendre en compte NLB sont `connection_time` et `tls_handshake_time`.

- Soyez prêt à interroger les journaux lorsque vous en aurez besoin. [Vous pouvez utiliser Amazon Athena pour interroger à la fois les ALB journaux et NLB les journaux.](#)
- Créez des alarmes pour les indicateurs liés aux performances, tels que [TargetResponseTime](#) pour ALB.

Ressources

Documents connexes :

- [ELB comparaison de produits](#)
- [AWS Infrastructure mondiale](#)
- [Amélioration des performances et réduction des coûts grâce à l'affinité des zones de disponibilité](#)
- [Procédure détaillée d'analyse des journaux avec Amazon Athena](#)
- [Interrogation des journaux de l'Application Load Balancer](#)
- [Surveillance de vos Application Load Balancers](#)
- [Surveillance de votre Network Load Balancer](#)
- [Utiliser Elastic Load Balancing pour répartir le trafic sur les instances dans votre groupe Auto Scaling](#)

Vidéos connexes :

- [AWS re:INVENT 2023 : Qu'est-ce que le réseau peut apporter à votre application ?](#)
- [AWS Re:inForce 20 : Comment utiliser Elastic Load Balancing pour améliorer votre niveau de sécurité à grande échelle](#)
- [AWS re:Invent 2018 : Elastic Load Balancing : analyse approfondie et meilleures pratiques](#)
- [AWS re:Invent 2021 - Comment choisir le bon équilibreur de charge pour vos charges de travail AWS](#)
- [AWS re:Invent 2019 : Tirez le meilleur parti d'Elastic Load Balancing pour différentes charges de travail](#)

Exemples connexes :

- [Équilibreur de charge de passerelle](#)
- [CDK et AWS CloudFormation des exemples pour l'analyse des journaux avec Amazon Athena](#)

PERF04-BP05 Choisissez les protocoles réseau pour améliorer les performances

Prenez des décisions concernant les protocoles de communication entre les systèmes et les réseaux en fonction de l'impact sur les performances de la charge de travail.

Il existe une relation entre la latence et la bande passante pour atteindre le débit. Si votre transfert de fichiers utilise le protocole de contrôle de transmission (TCP), des latences plus élevées réduiront probablement le débit global. Il existe des approches pour résoudre ce problème en TCP ajustant et en optimisant les protocoles de transfert, mais l'une des solutions consiste à utiliser le protocole User Datagram (UDP).

Anti-modèles courants :

- Vous l'utilisez TCP pour toutes les charges de travail, quelles que soient les exigences de performance.

Avantages liés au respect de cette bonne pratique : vérifiez que vous utilisez un protocole approprié pour la communication entre les utilisateurs et les composants de la charge de travail, afin d'améliorer l'expérience globale des utilisateurs de vos applications. Par exemple, le mode sans connexion UDP permet une vitesse élevée, mais il n'offre pas de retransmission ni une fiabilité élevée. TCP est un protocole complet, mais il nécessite une charge plus importante pour le traitement des paquets.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Si vous avez la possibilité de choisir différents protocoles pour votre application et que vous possédez l'expertise nécessaire dans ce domaine, optimisez votre application et l'expérience de l'utilisateur final en utilisant un autre protocole. Notez que cette approche présente des difficultés importantes et ne doit être tentée que si vous avez d'abord optimisé votre application à d'autres égards.

Pour améliorer les performances de votre charge de travail, il est essentiel de comprendre les exigences en matière de latence et de débit, puis de choisir des protocoles réseau qui optimisent les performances.

Quand envisager d'utiliser TCP

TCP fournit des données fiables et peut être utilisé pour la communication entre les composants de la charge de travail lorsque la fiabilité et la garantie de livraison des données sont importantes. De nombreuses applications Web s'appuient sur des protocoles TCP basés, tels que HTTP et HTTPS, pour ouvrir des TCP sockets permettant la communication entre les composants de l'application. Le transfert de données de courrier électronique et de fichiers est une application courante qui est également utilisée TCP, car il s'agit d'un mécanisme de transfert simple et fiable entre les composants de l'application. L'utilisation de TLS with TCP peut alourdir la communication, ce qui peut entraîner une augmentation de la latence et une réduction du débit, mais elle présente l'avantage de la sécurité. La surcharge provient principalement de la charge supplémentaire du processus de liaison, qui peut prendre plusieurs allers-retours pour se terminer. Une fois la liaison établie, la charge de chiffrement et de déchiffrement des données devient relativement faible.

Quand envisager d'utiliser UDP

UDP est un connection-less-oriented protocole et convient donc aux applications nécessitant une transmission rapide et efficace, telles que les données de journalisation, de surveillance et de VoIP. Pensez également à les utiliser UDP si vous disposez de composants de charge de travail qui répondent à de petites requêtes provenant d'un grand nombre de clients afin de garantir des performances optimales de la charge de travail. Datagram Transport Layer Security (DTLS) est l'UDP équivalent de Transport Layer Security (TLS). Lors de l'utilisation DTLS avec UDP, la surcharge provient du chiffrement et du déchiffrement des données, car le processus de prise de contact est simplifié. DTLS ajoute également une petite surcharge aux UDP paquets, car il inclut des champs supplémentaires pour indiquer les paramètres de sécurité et détecter les altérations.

Quand envisager d'utiliser SRD

Le datagramme fiable évolutif (SRD) est un protocole de transport réseau optimisé pour les charges de travail à haut débit en raison de sa capacité à équilibrer la charge du trafic sur plusieurs chemins et à récupérer rapidement en cas de perte de paquets ou de défaillance de liaison. SRD est donc mieux utilisé pour les charges de travail de calcul haute performance (HPC) qui nécessitent une communication à haut débit et à faible latence entre les nœuds de calcul. Il peut s'agir de tâches de traitement parallèle telles que la simulation, la modélisation et l'analyse de données qui impliquent un transfert important de données entre les nœuds.

Étapes d'implémentation

- Utilisez les services [AWS Global Accelerator](#) et [AWS Transfer Family](#) pour améliorer le débit de vos applications de transfert de fichiers en ligne. Le AWS Global Accelerator service vous aide à réduire le temps de latence entre vos appareils clients et votre charge de travail AWS. Vous

pouvez utiliser des AWS Transfer Family protocoles TCP basés tels que Secure Shell File Transfer Protocol (SFTP) et File Transfer Protocol over SSL (FTPS) pour dimensionner et gérer en toute sécurité vos transferts de fichiers vers les services AWS de stockage.

- Utilisez la latence du réseau pour déterminer si la communication entre les composants de la charge de travail TCP est appropriée. Si la latence du réseau entre votre application cliente et votre serveur est élevée, la prise de TCP contact à trois peut prendre un certain temps, ce qui a un impact sur la réactivité de votre application. Des mesures telles que le délai jusqu'au premier octet (TTFB) et le temps d'aller-retour (RTT) peuvent être utilisées pour mesurer la latence du réseau. Si votre charge de travail fournit du contenu dynamique aux utilisateurs, pensez à utiliser [Amazon CloudFront](#), qui établit une connexion permanente à chaque origine pour le contenu dynamique afin de supprimer le temps de configuration de la connexion qui ralentirait autrement chaque demande du client.
- L'utilisation TLS avec TCP ou UDP peut entraîner une augmentation de la latence et une réduction du débit pour votre charge de travail en raison de l'impact du chiffrement et du déchiffrement. Pour de telles charges de travail, pensez à SSL TLS /offloading sur [Elastic Load Balancing](#) afin d'améliorer les performances de la charge de travail en permettant à l'équilibreur de charge de gérer SSL les processus de TLS chiffrement et de déchiffrement au lieu de laisser les instances principales s'en charger. Cela peut contribuer à réduire l'utilisation des instances principales, ce qui peut améliorer les performances et augmenter la capacité.
- Utilisez le [Network Load Balancer \(NLB\)](#) pour déployer des services qui s'appuient sur le UDP protocole, tels que l'authentification et l'autorisation, la journalisation, l'DNSIoT et le streaming multimédia, afin d'améliorer les performances et la fiabilité de votre charge de travail. Le UDP trafic NLB entrant est réparti sur plusieurs cibles, ce qui vous permet d'adapter votre charge de travail horizontalement, d'augmenter la capacité et de réduire les frais généraux d'une seule cible.
- Pour vos charges de travail informatiques à hautes performances (HPC), pensez à utiliser la fonctionnalité [Elastic Network Adapter \(ENA\) Express](#) qui utilise le SRD protocole pour améliorer les performances du réseau en fournissant une bande passante à flux unique plus élevée (25 Gbit/s) et une latence de queue plus faible (99,9 centile) pour le trafic réseau entre les instances. EC2
- Utilisez l'[Application Load Balancer \(ALB\)](#) pour acheminer et équilibrer la charge de votre trafic gRPC (Remote Procedure Calls) entre les composants de la charge de travail ou entre les RPC clients et les services gRPC. gRPC utilise le protocole de transport TCP basé sur HTTP/2 et offre des avantages en termes de performances tels qu'un encombrement réseau réduit, une compression, une sérialisation binaire efficace, la prise en charge de nombreuses langues et un streaming bidirectionnel.

Ressources

Documents connexes :

- [Comment acheminer le UDP trafic vers Kubernetes](#)
- [Application Load Balancer](#)
- [EC2 Mise en réseau améliorée sous Linux](#)
- [EC2 Mise en réseau améliorée sous Windows](#)
- [EC2 Groupes de placement](#)
- [Activation de la mise en réseau améliorée avec l'adaptateur réseau Elastic \(ENA\) sur les instances Linux](#)
- [Network Load Balancer](#)
- [Produits de mise en réseau avec AWS](#)
- [Transition vers le routage basé sur la latence dans Amazon Route 53](#)
- [VPC Points de terminaison](#)

Vidéos connexes :

- [AWS re:Invent 2022 — Augmenter les performances du réseau sur les instances Amazon Elastic Compute Cloud de nouvelle génération](#)
- [AWS re:Invent 2022 — Les bases de la mise en réseau des applications](#)

Exemples connexes :

- [AWS Transit Gateway et solutions de sécurité évolutives](#)
- [Ateliers sur la mise en réseau AWS](#)

PERF04-BP06 Choisissez l'emplacement de votre charge de travail en fonction des exigences du réseau

Évaluez les options de placement des ressources afin de réduire la latence du réseau et d'améliorer le débit, offrant ainsi une expérience utilisateur optimale en réduisant les temps de chargement des pages et de transfert des données.

Anti-modèles courants :

- Vous regroupez toutes les ressources de charge de travail dans un seul emplacement géographique.
- Vous avez choisi la région la plus proche de votre emplacement, pas celle de l'utilisateur final de la charge de travail.

Avantages liés au respect de cette bonne pratique : l'expérience utilisateur est fortement affectée par le temps de latence entre l'utilisateur et votre application. En utilisant un réseau mondial AWS privé Régions AWS et approprié, vous pouvez réduire le temps de latence et offrir une meilleure expérience aux utilisateurs distants.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Les ressources, telles que EC2 les instances Amazon, sont placées dans des zones de disponibilité [Régions AWS](#) internes [AWS Outposts](#), [des zones AWS locales](#) ou [AWS Wavelength](#) des zones. Le choix de cet emplacement influence la latence et le débit du réseau à partir d'un emplacement donné de l'utilisateur. Les services périphériques tels qu'[Amazon CloudFront AWS Global Accelerator](#) peuvent également être utilisés pour améliorer les performances du réseau soit en mettant en cache le contenu sur des sites périphériques, soit en fournissant aux utilisateurs un chemin optimal vers la charge de travail via le réseau AWS mondial.

Amazon EC2 propose des groupes de placement pour la mise en réseau. Un groupe de placement est un regroupement logique d'instances permettant de réduire la latence. L'utilisation de groupes de placement dotés de types d'instances compatibles et d'un adaptateur réseau élastique (ENA) permet aux charges de travail de participer à un réseau de 25 Gbit/s à faible latence et à instabilité réduite. Les groupes de placement sont recommandés pour les charges de travail nécessitant une latence réseau faible, un débit réseau élevé ou les deux.

[Les services sensibles à la latence sont fournis sur des sites périphériques via un réseau AWS mondial, tel qu'Amazon. CloudFront](#) Ces emplacements périphériques fournissent généralement des services tels que le réseau de diffusion de contenu (CDN) et le système de noms de domaine (DNS). En disposant de ces services à la périphérie, les charges de travail peuvent répondre avec une faible latence aux demandes de contenu ou de DNS résolution. Ces services fournissent également des services géographiques tels que le ciblage géographique du contenu (qui fournit des contenus

différents en fonction de l'emplacement des utilisateurs finaux) ou le routage en fonction de la latence pour diriger les utilisateurs finaux vers la région plus proche (latence minimum).

Utilisez des services en périphérie pour réduire la latence et permettre la mise en cache de contenu. Configurez correctement le contrôle du cache pour les deux DNS et HTTP/HTTPS afin de tirer le meilleur parti de ces approches.

Étapes d'implémentation

- Capturez des informations sur le trafic IP entrant et sortant des interfaces réseau.
 - [Enregistrement du trafic IP à l'aide de VPC Flow Logs](#)
 - [Comment l'adresse IP du client est-elle préservée dans AWS Global Accelerator](#)
- Analysez les modèles d'accès au réseau dans votre charge de travail afin d'identifier comment les utilisateurs utilisent votre application.
 - Utilisez des outils de surveillance, tels qu'[Amazon CloudWatch](#) [AWS CloudTrail](#), pour recueillir des données sur les activités du réseau.
 - Analysez les données pour identifier le modèle d'accès au réseau.
- Choisissez les régions pour le déploiement de votre charge de travail en fonction des éléments clés suivants :
 - Lieu de stockage de vos données : pour les applications utilisant de grandes quantités de données (telles que le big data et le machine learning). Le code de l'application doit s'exécuter aussi près que possible des données.
 - Lieu de stockage de vos données : pour les applications orientées utilisateur, choisissez une région (ou des régions) proche des utilisateurs de votre charge de travail.
 - Autres contraintes : tenez compte des contraintes telles que le coût et la conformité, comme expliqué dans la section [Éléments à prendre en compte lors de la sélection d'une région pour vos charges de travail](#).
- Utilisez des zones locales [AWS](#) pour exécuter des charges de travail telles que le rendu vidéo. Les zones locales vous permettent de profiter des avantages liés à la présence de ressources de calcul et de stockage plus proches des utilisateurs finaux.
- Utilisez [AWS Outposts](#) pour les charges de travail qui doivent rester sur site et dont vous souhaitez qu'elles fonctionnent de manière transparente avec le reste de vos charges de travail dans AWS.
- Les applications telles que le streaming vidéo en direct haute résolution, le son haute fidélité et la réalité augmentée ou virtuelle (AR/VR) nécessitent ultra-low-latency des appareils 5G. Pour de telles applications, considérez [AWS Wavelength](#). AWS Wavelength intègre des services de AWS

calcul et de stockage dans les réseaux 5G, fournissant une infrastructure informatique de pointe mobile pour le développement, le déploiement et la mise à l'échelle d' ultra-low-latency applications.

- Utilisez des solutions de mise en cache locale ou [proposées par AWS](#) pour les ressources fréquemment utilisées afin d'améliorer les performances, de réduire les déplacements de données et de diminuer l'impact environnemental.

Service	Utilisation
Amazon CloudFront	Utilisez-le pour mettre en cache du contenu statique tel que des images, des scripts et des vidéos, ainsi que du contenu dynamique tel que API des réponses ou des applications Web.
Amazon ElastiCache	Permet de mettre en cache du contenu pour les applications Web.
DynamoDB Accelerator	Permet d'ajouter une accélération en mémoire à vos tables DynamoDB.

- Utilisez des services capables de vous aider à exécuter le code plus près des utilisateurs de votre charge de travail, tels que les suivants :

Service	Utilisation
Lambda@Edge	Destiné aux opérations exigeantes en puissance de calcul qui sont lancées lorsque des objets ne sont pas dans le cache.
CloudFront Fonctions Amazon	À utiliser pour des cas d'utilisation simples tels que HTTP des requêtes ou des manipulations de réponses qui peuvent être initiées par des fonctions de courte durée.
AWS IoT Greengrass	Permet d'exécuter du calcul local, une messagerie et une mise en cache de données pour les appareils connectés.

- Certaines applications nécessitent des points d'entrée fixes ou des performances plus élevées en réduisant la latence et l'instabilité du premier octet et en augmentant le débit. Ces applications peuvent bénéficier de services réseau qui fournissent des adresses IP anycast statiques et des TCP terminaisons aux emplacements périphériques. [AWS Global Accelerator](#) peut améliorer les performances de vos applications jusqu'à 60 % et permettre un basculement rapide pour les architectures multirégionales. AWS Global Accelerator vous fournit des adresses IP anycast statiques qui servent de point d'entrée fixe pour vos applications hébergées dans une ou plusieurs d' Régions AWS entre elles. Ces adresses IP permettent au trafic de pénétrer sur le réseau AWS mondial aussi près que possible de vos utilisateurs. AWS Global Accelerator réduit le temps de configuration de la connexion initiale en établissant une TCP connexion entre le client et l'emplacement AWS périphérique le plus proche du client. Passez en revue l'utilisation de AWS Global Accelerator pour améliorer les performances de vos TCP/UDP workloads et permettre un basculement rapide pour les architectures multirégionales.

Ressources

Bonnes pratiques associées :

- [COST07-BP02 Mettre en œuvre les régions en fonction des coûts](#)
- [COST08-BP03 Mettre en œuvre des services pour réduire les coûts de transfert de données](#)
- [REL10-BP01 Déployer la charge de travail sur plusieurs sites](#)
- [REL10-BP02 Sélectionnez les emplacements appropriés pour votre déploiement multisite](#)
- [SUS01-BP01 Choisissez la région en fonction des exigences commerciales et des objectifs de durabilité](#)
- [SUS02-BP04 Optimiser le placement géographique des charges de travail en fonction de leurs exigences en matière de réseau](#)
- [SUS04-BP07 Minimiser le mouvement des données sur les réseaux](#)

Documents connexes :

- [AWS Infrastructure mondiale](#)
- [AWS Zones locales et AWS Outposts choix de la technologie adaptée à votre charge de travail périphérique](#)
- [Groupes de placement](#)
- [AWS Zones Locales](#)

- [AWS Outposts](#)
- [AWS Wavelength](#)
- [Amazon CloudFront](#)
- [AWS Global Accelerator](#)
- [AWS Direct Connect](#)
- [AWS Site-to-Site VPN](#)
- [Amazon Route 53](#)

Vidéos connexes :

- [AWS Vidéo explicative sur les Zones Locales](#)
- [AWS Outposts: Overview and How it Works](#)
- [AWS re:Invent 2023 - Une stratégie de migration pour les charges de travail en périphérie et sur site](#)
- [AWS re:INVENT 2021 - AWS Outposts : Apporter l' AWS expérience sur site](#)
- [AWS re:Invent 2020 : AWS Wavelength : Exécutez des applications avec une latence extrêmement faible à la périphérie de la 5G](#)
- [AWS re:Invent 2022 - Zones AWS locales : création d'applications pour une périphérie distribuée](#)
- [AWS re:Invent 2021 - Création de sites Web à faible latence avec Amazon CloudFront](#)
- [AWS re:Invent 2022 - Améliorez les performances et la disponibilité avec AWS Global Accelerator](#)
- [AWS re:Invent 2022 - Construisez votre réseau étendu mondial en utilisant AWS](#)
- [AWS re:Invent 2020 : gestion du trafic mondial avec Amazon Route 53](#)

Exemples connexes :

- [AWS Global Accelerator Atelier de routage personnalisé](#)
- [Gestion des réécritures et des redirections à l'aide des fonctions de périphérie](#)

PERF04-BP07 Optimiser la configuration du réseau en fonction des métriques

Utilisez les données collectées et analysées pour prendre des décisions avisées concernant l'optimisation de votre configuration réseau.

Anti-modèles courants :

- Vous supposez que tous les problèmes liés aux performances sont liés à l'application.
- Vous testez uniquement les performances de votre réseau à partir d'un emplacement proche de l'endroit où vous avez déployé la charge de travail.
- Vous utilisez des configurations par défaut pour tous les services du réseau.
- Vous surdimensionnez la ressource réseau afin de fournir une capacité suffisante.

Avantages liés au respect de cette bonne pratique : la collecte des métriques nécessaires de votre réseau AWS et la mise en œuvre d'outils de surveillance du réseau vous permettent de comprendre les performances du réseau et d'optimiser les configurations du réseau.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : faible

Directives d'implémentation

La surveillance du trafic en provenance VPCs et à destination des sous-réseaux ou des interfaces réseau est essentielle pour comprendre comment utiliser les ressources AWS réseau et optimiser les configurations réseau. À l'aide des outils AWS réseau suivants, vous pouvez examiner plus en détail les informations relatives à l'utilisation du trafic, à l'accès au réseau et aux journaux.

Étapes d'implémentation

- Identifiez les indicateurs de performance clés tels que la latence ou la perte de paquets à collecter. AWS fournit plusieurs outils qui peuvent vous aider à collecter ces statistiques. Les outils suivants vous permettent d'obtenir des informations supplémentaires sur l'utilisation du trafic, l'accès au réseau et les journaux.

AWS outil	Où utiliser
Gestionnaire d'adresses VPC IP Amazon.	IPAM Utilisez-le pour planifier, suivre et surveiller les adresses IP pour vos charges de travail AWS et celles sur site. Il s'agit d'une bonne pratique pour optimiser l'utilisation et l'allocation des adresses IP.
VPC Journaux de flux	Utilisez les journaux de VPC flux pour capturer des informations détaillées sur le trafic à destination et en provenance des interfaces réseau de votre VPCs. Avec VPC Flow Logs, vous pouvez diagnostiquer les règles de groupe de sécurité trop restrictives ou trop permissives et déterminer la direction du trafic à destination et en provenance des interfaces réseau.
AWS Transit Gateway Journaux de flux	Utilisez les journaux de AWS Transit Gateway flux pour capturer des informations sur le trafic IP à destination et en provenance de vos passerelles de transit.
DNS journalisation des requêtes	Enregistrez les informations relatives aux DNS requêtes publiques ou privées reçues par Route 53. Les journaux vous permettent d'optimiser les DNS configurations en comprenant le domaine ou le sous-domaine qui a été demandé ou les EDGE emplacements Route 53 qui ont répondu aux DNS requêtes.

AWS outil	Où utiliser
Reachability Analyzer	<p>Reachability Analyzer vous aide à analyser et à déboguer l'accessibilité du réseau. Reachability Analyzer est un outil d'analyse de configuration qui vous permet d'effectuer des tests de connectivité entre une ressource source et une ressource de destination dans votre VPC. Cet outil vous aide à vérifier que votre configuration réseau correspond à la connectivité souhaitée.</p>
Analyseur d'accès réseau	<p>Vous pouvez utiliser l'Analyseur d'accès réseau pour comprendre l'accès réseau à vos ressources. Vous pouvez utiliser l'analyseur d'accès réseau pour spécifier vos exigences en matière d'accès au réseau et identifier les chemins d'accès potentiels qui ne répondent pas à vos exigences spécifiées. En optimisant la configuration de votre réseau correspondant, vous pouvez comprendre et vérifier l'état de votre réseau et démontrer si votre réseau sur AWS répond à vos exigences de conformité.</p>

AWS outil	Où utiliser
Amazon CloudWatch	Utilisez Amazon CloudWatch et activez les métriques appropriées pour les options de réseau. Veillez à choisir la métrique de réseau adaptée à votre charge de travail. Par exemple, vous pouvez activer les métriques pour l'utilisation des adresses VPC réseau, la VPC NAT passerelle AWS Transit Gateway, le VPN tunnel AWS Network Firewall, Elastic Load Balancing et AWS Direct Connect. La surveillance continue des métriques est une bonne pratique pour observer et comprendre l'état et l'utilisation de votre réseau. Elle vous aide à optimiser la configuration du réseau en fonction de vos observations.
AWS Network Manager	Vous pouvez ainsi surveiller les performances historiques et en temps réel du réseau AWS mondial à des fins opérationnelles et de planification. AWS Network Manager Network Manager fournit une latence réseau globale entre les zones de disponibilité Régions AWS et au sein de chaque zone de disponibilité, ce qui vous permet de mieux comprendre le lien entre les performances de votre application et les performances du AWS réseau sous-jacent.
Amazon CloudWatch RUM	Utilisez Amazon CloudWatch RUM pour collecter les statistiques qui vous fournissent les informations qui vous aideront à identifier, à comprendre et à améliorer l'expérience utilisateur.

- Identifiez les principaux intervenants et les modèles de trafic des applications à l'aide VPC des journaux de AWS Transit Gateway flux.

- Évaluez et optimisez votre architecture réseau actuelle VPCs, y compris les sous-réseaux et le routage. Par exemple, vous pouvez évaluer dans quelle mesure le VPC peering est différent ou vous AWS Transit Gateway aider à améliorer la mise en réseau dans votre architecture.
- Évaluez les chemins de routage de votre réseau pour vérifier que le chemin le plus court entre les destinations est toujours utilisé. L'Analyseur d'accès réseau vous aide à le faire.

Ressources

Documents connexes :

- [Journalisation des DNS requêtes publiques](#)
- [Qu'est-ce que c'est IPAM ?](#)
- [Définir Reachability Analyzer](#)
- [Définir l'Analyseur d'accès réseau](#)
- [CloudWatch indicateurs pour votre VPCs](#)
- [Optimisez les performances et réduisez les coûts d'analyse du réseau avec VPC Flow Logs au format Apache Parquet](#)
- [Surveillance de vos réseaux mondiaux et principaux à l'aide des CloudWatch métriques Amazon](#)
- [Surveiller en permanence le trafic et les ressources du réseau](#)

Vidéos connexes :

- [AWS re:Invent 2023 — Guide du développeur sur les réseaux cloud](#)
- [AWS re:Invent 2023 — Prêts pour la suite ? Designing networks for growth and flexibility](#)
- [AWS re:Invent 2023 — VPC Designs avancés et nouvelles fonctionnalités](#)
- [AWS re:Invent 2022 — Approfondissez l'infrastructure réseau AWS](#)
- [AWS re:Invent 2020 — Meilleures pratiques et astuces de mise en réseau avec le cadre Well-Architected AWS](#)
- [AWS re:Invent 2020 — Surveillance et résolution des problèmes du trafic réseau](#)

Exemples connexes :

- [Ateliers sur la mise en réseau AWS](#)
- [Surveillance réseau AWS](#)

- [Observation et diagnostic de votre réseau sur AWS](#)
- [Détecter et corriger les erreurs de configuration du réseau sur AWS](#)

Processus et culture

Lors de la création de l'architecture des charges de travail, vous pouvez adopter certains principes et certaines pratiques pour optimiser l'exécution de charges de travail cloud efficaces et performantes. Ce domaine d'intérêt propose les bonnes pratiques pour l'adoption d'une culture qui favorise l'efficacité des performances des charges de travail dans le cloud.

Tenez compte de ces principes clés pour développer cette culture :

- **Infrastructure en tant que code** : définissez votre infrastructure en tant que code à l'aide de méthodes telles que les modèles AWS CloudFormation. L'utilisation de modèles vous permet de placer votre infrastructure en mode de contrôle de code source parallèlement au code et aux configurations de votre application. Vous pouvez ainsi appliquer les pratiques utilisées pour développer des logiciels à votre infrastructure et itérer rapidement.
- **Pipeline de déploiement** : utilisez un pipeline de déploiement d'intégration continue (CI) et de livraison continue (CD) (par exemple, référentiel de code source, systèmes de génération, déploiement et automatisation des tests) pour déployer votre infrastructure. Vous pouvez ainsi déployer de manière reproductible et cohérente, le tout à un faible coût, à mesure que vous itérez.
- **Métriques bien définies** : configurez vos métriques et votre solution de surveillance pour capturer les indicateurs de performances clés (KPI). Nous vous recommandons d'utiliser des métriques techniques, mais aussi des métriques commerciales. Pour les sites Web ou les applications mobiles, les indicateurs clés capturent le temps jusqu'au premier octet ou le rendu. D'autres mesures généralement applicables comprennent le nombre de threads, le taux de récupérateur de mémoire et les états d'attente. Les métriques commerciales, telles que les coûts cumulés agrégés par demande, peuvent vous permettre d'identifier des solutions pour réduire vos coûts. Réfléchissez bien à la façon dont vous prévoyez d'interpréter les métriques. Par exemple, vous pouvez choisir le maximum ou le 99e centile plutôt que la moyenne.
- **Tests de performance automatiques** : dans le cadre de votre processus de déploiement, des tests de performance peuvent se déclencher automatiquement une fois les tests en cours d'exécution bien effectués. L'automatisation doit créer un environnement, configurer des conditions initiales (comme des données de test), puis exécuter une série d'analyses comparatives et de tests de charge. Les résultats de ces tests doivent être rattachés à la version de génération afin que vous puissiez suivre l'évolution des performances dans le temps. Pour les tests de longue durée, vous pouvez rendre cette partie du pipeline asynchrone par rapport au reste de la compilation. Sinon, vous pouvez exécuter des tests de performances pendant la nuit en utilisant les instances Spot Amazon EC2.

- **Génération de charge** : vous devez créer une série de scripts qui reproduisent des parcours utilisateur synthétiques ou préenregistrés. Ces scripts doivent être idempotents et non couplés. Il se peut que vous deviez aussi inclure à cette série des scripts de préparation pour obtenir des résultats valides. Dans la mesure du possible, vos scripts de test doivent pouvoir répliquer le comportement d'utilisation en production. Vous pouvez utiliser un logiciel ou des solutions de logiciel en tant que service (SaaS) pour générer la charge. Envisagez d'utiliser les solutions [AWS Marketplace](#) et les [instances Spot](#) : elles peuvent être des moyens économiques de générer la charge.
- **Visibilité des performances** : les métriques clés doivent être visibles pour votre équipe, en particulier pour chaque version. Vous pouvez ainsi identifier les tendances positives ou négatives significatives au fil du temps. Vous devez également afficher les métriques sur le nombre d'erreurs ou d'exceptions pour vous assurer que vous testez un système fonctionnel.
- **Visualisation** : utilisez des techniques de visualisation qui permettent d'identifier clairement l'origine des problèmes de performances, les points chauds, les états d'attente ou les taux d'utilisation faibles. Superposez les métriques de performance sur les schémas d'architecture, des graphiques ou codes d'appel qui peuvent vous aider à identifier rapidement les problèmes.
- **Processus d'examen régulier** : les architectures qui présentent des performances médiocres sont généralement le résultat d'un processus d'évaluation des performances inexistant ou interrompu. Si votre architecture est peu performante, la mise en œuvre d'un processus d'évaluation des performances vous permet de procéder à des améliorations itératives.
- **Optimisation continue** : adoptez une culture permettant d'optimiser en permanence l'efficacité des performances de votre charge de travail dans le cloud.

Bonnes pratiques

- [PERF05-BP01 Définition d'indicateurs de rendement clés \(KPI\) pour mesurer l'état et les performances de la charge de travail](#)
- [PERF05-BP02 Utiliser des solutions de surveillance pour comprendre les domaines dans lesquels les performances sont les plus critiques](#)
- [PERF05-BP03 Définir un processus pour améliorer les performances de la charge de travail](#)
- [PERF05-BP04 Testez votre charge de travail](#)
- [PERF05-BP05 Utiliser l'automatisation pour résoudre de manière proactive les problèmes liés aux performances](#)
- [PERF05-BP06 Maintenez votre charge de travail et vos services up-to-date](#)
- [PERF05-BP07 Vérification des métriques à intervalles réguliers](#)

PERF05-BP01 Définition d'indicateurs de rendement clés (KPI) pour mesurer l'état et les performances de la charge de travail

Identifiez les KPI qui mesurent les performances de la charge de travail de manière quantitative et qualitative. Les KPI vous aident à mesurer l'état et les performances d'une charge de travail par rapport à un objectif métier.

Anti-modèles courants :

- Vous surveillez uniquement les métriques au niveau du système pour avoir un aperçu de votre charge de travail et ne comprenez pas les impacts commerciaux de ces métriques.
- Vous supposez que vos KPI sont déjà en cours de publication et de partage en tant que données de métriques standard.
- Vous ne définissez pas de KPI quantitatif et mesurable.
- Vous ne tenez pas compte des objectifs ni des stratégies de l'entreprise pour définir vos KPI.

Avantages liés au respect de cette bonne pratique : en identifiant les KPI spécifiques qui représentent l'état et les performances de la charge de travail, vous pouvez aligner les équipes sur leurs priorités et définir des résultats commerciaux atteignables. Le partage de ces métriques avec tous les départements offre une visibilité et un alignement sur les seuils, les attentes et l'impact commercial.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

Les KPI permettent aux équipes commerciales et d'ingénierie de s'aligner sur la mesure des objectifs et des stratégies et sur la façon dont ces facteurs se combinent pour générer des résultats commerciaux. Par exemple, une charge de travail de site Web peut utiliser le temps de chargement de la page comme indication des performances globales. Cette métrique serait l'un des éléments de données pris en compte qui mesure l'expérience d'un utilisateur. En plus d'identifier les temps limites de chargement des pages, vous devez documenter le résultat attendu ou le risque commercial si les performances idéales ne sont pas atteintes. Un temps de chargement long des pages affecte directement vos utilisateurs finaux, nuit à leur expérience utilisateur et peut entraîner une perte de clients. Lorsque vous définissez vos seuils de KPI, combinez à la fois les points de référence en vigueur dans votre secteur et les attentes de vos utilisateurs finaux. Par exemple, si le point de référence actuel établi par votre secteur d'activité pour le chargement d'une page Web est un délai de

deux secondes, mais que vos utilisateurs finaux s'attendent à ce qu'une page Web se charge dans un délai d'une seconde, vous devez prendre en compte ces deux éléments de données lors de la définition des KPI.

Votre équipe doit évaluer les KPI de votre charge de travail à l'aide de données précises en temps réel et de données historiques à titre de référence et créer des tableaux de bord qui effectuent des calculs de métriques par rapport à vos données de KPI pour générer des informations opérationnelles et d'utilisation. Les KPI doivent être documentés et inclure les seuils qui soutiennent les objectifs et les stratégies de l'entreprise et doivent être mappés aux métriques surveillées. Les KPI doivent être revus lorsque les objectifs commerciaux, les stratégies ou les exigences des utilisateurs finaux changent.

Étapes d'implémentation

- Identification des parties prenantes : identifiez et documentez les principales parties prenantes de l'entreprise, y compris les équipes de développement et d'exploitation.
- Définition d'objectifs : collaborez avec ces parties prenantes pour définir et documenter les objectifs de votre charge de travail. Tenez compte des aspects critiques des performances de vos charges de travail, tels que le débit, le temps de réponse et le coût, ainsi que des objectifs métier, tels que la satisfaction des utilisateurs.
- Passage en revue des bonnes pratiques du secteur : passez en revue les bonnes pratiques du secteur pour identifier les KPI pertinents qui correspondent à vos objectifs en matière de charge de travail.
- Identification des métriques : identifiez les métriques qui correspondent aux objectifs de votre charge de travail et qui peuvent vous aider à mesurer les performances et les objectifs commerciaux. Établissez des KPI sur la base de ces métriques. Les mesures telles que le temps de réponse moyen ou le nombre d'utilisateurs simultanés sont des exemples de métriques.
- Définition et documentation des KPI : utilisez les bonnes pratiques du secteur et les objectifs de votre charge de travail pour définir des cibles pour votre KPI de charge de travail. Utilisez ces informations pour définir les seuils de KPI pour les niveaux de gravité ou d'alarme. Identifiez et documentez le risque et l'impact du non-respect d'un KPI.
- Mise en œuvre de la surveillance : utilisez des outils de surveillance tels qu'[Amazon CloudWatch](#) ou [AWS Config](#) pour collecter des métriques et mesurer les KPI.
- Communication visuelle des KPI : utilisez des outils de tableau de bord tels qu'[Amazon QuickSight](#) pour visualiser et communiquer les indicateurs de performance clés aux parties prenantes.

- **Analyse et optimisation** : passez en revue et analysez régulièrement les KPI pour identifier les domaines de votre charge de travail qui doivent être améliorés. Collaborez avec les parties prenantes pour mettre en œuvre ces améliorations.
- **Révision et affinage** : passez régulièrement en revue les indicateurs et les indicateurs de performance clés pour évaluer leur efficacité, en particulier lorsque les objectifs commerciaux ou les performances de la charge de travail changent.

Ressources

Documents connexes :

- [Documentation CloudWatch](#)
- [Surveillance, journalisation et performances AWS Partner](#)
- [Outils d'observabilité d'AWS](#)
- [L'importance des indicateurs de rendement clés \(KPI\) pour les migrations vers le cloud à grande échelle](#)
- [Suivi des KPI d'optimisation des coûts avec KPI Dashboard](#)
- [Documentation X-Ray](#)
- [Utilisation des tableaux de bord Amazon CloudWatch](#)
- [KPI QuickSight](#)

Vidéos connexes :

- [AWS re:Invent 2023 - Optimize cost and performance and track progress toward mitigation](#)
- [AWS re:Invent 2023 - Manage resource lifecycle events at scale with AWS Health](#)
- [AWS re:Invent 2023 - Performances & efficiency at Pinterest: Optimizing the latest instances](#)
- [AWS re:Invent 2022 - AWS optimization: Actionable steps for immediate results](#)
- [AWS re:Invent 2023 – Building an effective observability strategy](#)
- [AWS Summit SF 2022 - Full-stack observability and application monitoring with AWS](#)
- [AWS re:Invent 2023 - Scaling on AWS for the first 10 million users](#)
- [AWS re:Invent 2022 - How Amazon uses better metrics for improved website performance](#)
- [Creating an Effective Metrics Strategy for Your Business | AWS Events](#)

Exemples connexes :

- [Création d'un tableau de bord avec QuickSight](#)

PERF05-BP02 Utiliser des solutions de surveillance pour comprendre les domaines dans lesquels les performances sont les plus critiques

Comprenez et identifiez les domaines où l'augmentation des performances de votre charge de travail aura un impact positif sur l'efficacité ou l'expérience client. Par exemple, un site Web qui comporte un grand nombre d'interactions clients pourrait gagner à utiliser des services de périphérie pour rapprocher la diffusion de contenus des clients.

Anti-modèles courants :

- Vous supposez que les mesures de calcul standard telles que CPU l'utilisation ou la pression de la mémoire sont suffisantes pour détecter les problèmes de performances.
- Vous n'utilisez que les métriques par défaut enregistrées par le logiciel de surveillance que vous avez sélectionné.
- Vous n'examinez les métriques qu'en cas de problème.

Avantages de l'établissement de cette meilleure pratique : la compréhension des domaines de performance critiques aide les responsables de la charge de travail à surveiller KPIs et à hiérarchiser les améliorations à fort impact.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : élevé

Directives d'implémentation

Configurez le end-to-end suivi pour identifier les modèles de trafic, la latence et les domaines de performance critiques. Surveillez vos modèles d'accès aux données afin d'identifier les requêtes lentes ou les données mal fragmentées et partitionnées. Identifiez les zones de charge de travail limitées à l'aide de tests ou de surveillance des charges.

améliorer l'efficacité des performances en comprenant votre architecture, vos modèles de trafic et d'accès aux données, et identifier vos temps de latence et de traitement. Identifier les goulots

d'étranglement potentiels qui pourraient avoir une incidence sur l'expérience client à mesure que la charge de travail augmente. Après avoir enquêté sur ces domaines, déterminez quelle solution vous pouvez déployer afin de surmonter ces problèmes de performances.

Étapes d'implémentation

- Configurez end-to-end la surveillance pour capturer tous les composants et mesures de la charge de travail. Voici des exemples de solutions de surveillance sur AWS.

Service	Où utiliser
Surveillance CloudWatch des utilisateurs réels d'Amazon () RUM	Pour capturer les métriques de performances des applications à partir de sessions réelles côté client et front-end.
AWS X-Ray	Pour tracer le trafic à travers les couches applicatives et identifier la latence entre les composants et les dépendances. Utilisez les cartographies de services X-Ray afin de voir les relations et la latence entre les composants de la charge de travail.
Informations sur les performances d'Amazon Relational Database Service	Pour consulter les métriques de performances de la base de données et identifier les améliorations des performances.
Surveillance RDS améliorée d'Amazon	Pour consulter les métriques de performances du système d'exploitation de la base de données.
Amazon DevOps Guru	Pour détecter les modèles de fonctionnement anormaux afin que vous puissiez identifier les problèmes opérationnels avant qu'ils n'affectent vos clients.

- Effectuez des tests afin de générer des métriques, d'identifier les tendances de trafic, les goulots d'étranglement et les domaines de performance critiques. Voici quelques exemples de méthodes de test :

- Configurez [CloudWatchSynthetic Canaries](#) pour imiter les activités des utilisateurs basées sur le navigateur de manière programmatique à l'aide de tâches cron Linux ou d'expressions de taux afin de générer des métriques cohérentes au fil du temps.
- Utiliser le [test de charge distribuéAWS](#) afin de générer un trafic de pointe ou de tester la charge de travail au taux de croissance attendu.
- Évaluez les métriques et la télémétrie pour identifier vos domaines de performances critiques. Examinez ces domaines avec votre équipe afin de discuter de la surveillance et des solutions pour éviter les goulots d'étranglement.
- Expérimentez des améliorations des performances et mesurez ces changements avec des données. Par exemple, vous pouvez utiliser [CloudWatchEvidently](#) pour tester les nouvelles améliorations et les impacts sur les performances de votre charge de travail.

Ressources

Documents connexes :

- [Quoi de neuf en matière d' AWS observabilité à re:Invent 2023](#)
- [Bibliothèque Amazon Builders' Library](#)
- [Documentation X-Ray](#)
- [Amazon CloudWatch RUM](#)
- [Amazon DevOps Guru](#)

Vidéos connexes :

- [AWS re:Invent 2023 - \[LAUNCH\] Surveillance des applications pour les charges de travail modernes](#)
- [AWS re:Invent 2023 - Mise en œuvre de l'observabilité des applications](#)
- [AWS re:Invent 2023 - Élaboration d'une stratégie d'observabilité efficace](#)
- [AWS Summit SF 2022 - Observabilité complète et surveillance des applications avec AWS](#)
- [AWS re:Invent 2022 - AWS optimisation : étapes réalisables pour des résultats immédiats](#)
- [AWS re:Invent 2022 - La bibliothèque Amazon Builders' Library : 25 ans d'excellence opérationnelle d'Amazon](#)
- [AWS re:Invent 2022 - Comment Amazon utilise de meilleurs indicateurs pour améliorer les performances de son site Web](#)

- [Surveillance visuelle des applications avec Amazon CloudWatch Synthetics](#)

Exemples connexes :

- [Mesurez le temps de chargement des pages avec Amazon CloudWatch Synthetics](#)
- [Client CloudWatch RUM Web Amazon](#)
- [X-Ray SDK pour Python](#)
- [Test de charge distribué sur AWS](#)

PERF05-BP03 Définir un processus pour améliorer les performances de la charge de travail

Définissez un processus d'évaluation de nouveaux services, modèles de conception, types de ressources et configurations au fur et à mesure qu'ils deviennent disponibles. Par exemple, exécutez des tests de performances existants sur de nouvelles offres d'instances afin de déterminer leur potentiel d'amélioration de votre charge de travail.

Anti-modèles courants :

- Vous supposez que votre architecture actuelle est statique et ne sera pas mise à jour au fil du temps.
- Vous introduisez des modifications d'architecture au fil du temps sans justification basée sur les métriques.

Avantages liés au respect de cette bonne pratique : un processus défini pour les modifications d'architecture rend possible l'utilisation des données collectées pour influencer la conception de votre charge de travail au fil du temps.

Niveau de risque encouru si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Les performances de votre charge de travail présentent quelques contraintes clés. Documentez-les pour connaître les types d'innovations qui pourraient améliorer les performances de votre charge de travail. Utilisez ces informations lors de l'apprentissage de nouveaux services ou la technologie au fur

et à mesure de leur disponibilité afin d'identifier les moyens d'atténuer des contraintes ou des goulets d'étranglement.

Identifiez les principales contraintes de performance pour votre charge de travail. Documentez les contraintes environnementales de votre charge de travail pour connaître les types d'innovations qui pourraient améliorer les performances de celle-ci.

Étapes d'implémentation

- Identifier KPIs : Identifiez les performances de votre charge de travail KPIs comme indiqué dans la section [PERF05-BP01 Définition d'indicateurs de rendement clés \(KPI\) pour mesurer l'état et les performances de la charge de travail](#) pour établir une base de référence de votre charge de travail.
- Mettre en œuvre le suivi : utilisez des [outils AWS d'observabilité](#) pour collecter des indicateurs de performance et les mesurer KPIs.
- Réalisation d'une analyse : effectuez une analyse approfondie pour identifier les domaines (tels que la configuration et le code d'application) de votre charge de travail qui ne sont pas performants, comme indiqué dans [PERF05-BP02 Utiliser des solutions de surveillance pour comprendre les domaines dans lesquels les performances sont les plus critiques](#). Utilisez vos outils d'analyse et de performance pour identifier les stratégies d'amélioration des performances.
- Validation des améliorations : utilisez des environnements de test (sandbox) ou en préproduction pour valider l'efficacité des stratégies d'amélioration.
- Mise en œuvre des modifications : mettez en œuvre les modifications en production et surveillez en permanence les performances de la charge de travail. Documentez les améliorations et communiquez-les aux parties prenantes.
- Révision et affinage : passez régulièrement en revue votre processus d'amélioration des performances afin d'identifier les domaines à améliorer.

Ressources

Documents connexes :

- [Blog AWS](#)
- [Quoi de neuf avec AWS](#)
- [AWS Générateur de compétences](#)

Vidéos connexes :

- [AWS re:Invent 2022 - Fournir des architectures durables et performantes](#)
- [AWS re:Invent 2023 - Optimisez les coûts et les performances et suivez les progrès en matière d'atténuation](#)
- [AWS re:Invent 2022 - AWS optimisation : étapes réalisables pour des résultats immédiats](#)
- [AWS re:Invent 2022 - Optimisez vos AWS charges de travail grâce à des conseils sur les meilleures pratiques](#)

Exemples connexes :

- [AWS Github](#)

PERF05-BP04 Testez votre charge de travail

Effectuez un test de charge de votre charge de travail pour vérifier qu'elle peut supporter la charge de production et identifier les éventuels goulots d'étranglement en termes de performances.

Anti-modèles courants :

- Vous testez les différentes parties et non la totalité de votre charge de travail.
- Vous testez la charge sur une infrastructure qui n'est pas la même que votre environnement de production.
- Vous n'effectuez le test de charge que pour la charge prévue sans aller au-delà, avec pour but de prévoir où vous pourriez rencontrer des problèmes à l'avenir.
- Vous effectuez des tests de charge sans consulter la [politique de EC2 test d'Amazon](#) et sans soumettre de formulaire de soumission d'événements simulés. Cela entraîne l'échec de votre test, car il ressemble à un denial-of-service événement.

Avantages liés au respect de cette bonne pratique : la mesure de vos performances dans le cadre d'un test de charge vous indiquera où vous serez affecté au fil de l'augmentation de la charge. Cela peut vous permettre d'anticiper les changements nécessaires avant qu'ils n'affectent votre charge de travail.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : faible

Directives d'implémentation

Les tests de charge dans le cloud sont un processus visant à mesurer les performances de la charge de travail cloud dans des conditions réalistes avec la charge utilisateur attendue. Ce processus implique la mise en service d'un environnement cloud de type production, l'utilisation d'outils de test de charge pour générer la charge et l'analyse de métriques pour évaluer la capacité de votre charge de travail à gérer une charge réaliste. Pour effectuer un test de charge, vous devez exécuter des versions de données de production factices ou légèrement altérées (supprimez les données sensibles ou les informations d'identification). Effectuez automatiquement des tests de charge dans le cadre de votre pipeline de livraison et comparez les résultats par rapport à des seuils KPIs et à des seuils prédéfinis. Ce processus vous permet de continuer à atteindre les performances requises.

Étapes d'implémentation

- Définition de vos objectifs de test : identifiez les aspects de performance de votre charge de travail que vous souhaitez évaluer, tels que le débit et le temps de réponse.
- Sélection d'un outil de test : choisissez et configurez l'outil de test de charge adapté à votre charge de travail.
- Configuration de votre environnement : configurez l'environnement de test en fonction de votre environnement de production. Vous pouvez utiliser AWS les services pour exécuter des environnements de production afin de tester votre architecture.
- Mettez en œuvre la surveillance : utilisez des outils de surveillance tels qu'[Amazon CloudWatch](#) pour collecter des métriques sur les ressources de votre architecture. Vous pouvez également collecter et publier des métriques personnalisées.
- Définition de des scénarios : définissez les scénarios et les paramètres de test de charge (tels que la durée du test et le nombre d'utilisateurs).
- Tests de charge : réalisez des scénarios de test à grande échelle. Profitez-en AWS Cloud pour tester votre charge de travail afin de découvrir où elle ne parvient pas à évoluer ou si elle évolue de manière non linéaire. Par exemple, utilisez les instances Spot pour générer des charges à faible coût et découvrir les goulots d'étranglement avant de les rencontrer en production.
- Analyse des résultats des tests : analysez les résultats pour identifier les goulots d'étranglement en matière de performances et les domaines à améliorer.
- Documentation et partage des résultats : documentez et rendez compte des résultats et des recommandations. Partagez ces informations avec les parties prenantes pour les aider à prendre des décisions éclairées concernant les stratégies d'optimisation des performances.

- Itération continue : les tests de charge doivent être effectués à une cadence régulière, en particulier après un changement ou une mise à jour du système.

Ressources

Documents connexes :

- [Amazon CloudWatch RUM](#)
- [Amazon CloudWatch Synthetics](#)
- [Test de charge distribué sur AWS](#)

Vidéos connexes :

- [AWS Sommet ANZ 2023 : Accélérez en toute confiance grâce AWS aux tests de charge distribués](#)
- [AWS re:Invent 2022 - Tirez parti AWS de vos 10 premiers millions d'utilisateurs](#)
- [Résoudre avec AWS des solutions : tests de charge distribués](#)
- [AWS re:Invent 2021 - Optimisez les applications grâce aux informations des utilisateurs finaux avec Amazon CloudWatch RUM](#)
- [Démonstration d'Amazon CloudWatch Synthetics](#)

Exemples connexes :

- [Test de charge distribué sur AWS](#)

PERF05-BP05 Utiliser l'automatisation pour résoudre de manière proactive les problèmes liés aux performances

Utilisez des indicateurs de performance clés (KPIs), combinés à des systèmes de surveillance et d'alerte, pour résoudre de manière proactive les problèmes liés aux performances.

Anti-modèles courants :

- Vous autorisez uniquement le personnel des opérations à apporter des modifications opérationnelles à la charge de travail.

- Vous confiez toutes les activités de filtre des alarmes à l'équipe des opérations sans correction proactive.

Avantages liés au respect de cette bonne pratique : la correction proactive des actions d'alarme permet au personnel d'assistance de se concentrer sur les éléments qui ne sont pas exploitables automatiquement. Cela permet au personnel des opérations de gérer toutes les alarmes sans être submergé et de se concentrer uniquement sur les alarmes critiques.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : faible

Directives d'implémentation

Utilisez des alarmes pour déclencher des actions automatisées afin de corriger les problèmes dans la mesure du possible. Faites remonter l'alarme aux personnes qui peuvent répondre si une réponse automatique n'est pas possible. Par exemple, vous pouvez disposer d'un système capable de prédire les valeurs attendues des indicateurs de performance clés (KPI) et de déclencher une alarme lorsqu'ils dépassent certains seuils, ou d'un outil capable d'arrêter ou d'annuler automatiquement les déploiements s'ils KPIs sont en dehors des valeurs attendues.

Mettez en place des processus qui rendent visibles les performances pendant que votre charge de travail est en cours d'exécution. Créez des tableaux de bord de surveillance et établissez des normes de référence pour les attentes en matière de performances pour déterminer si les performances de la charge de travail sont optimales.

Étapes d'implémentation

- Identification du processus de remédiation : identifiez et comprenez le problème lié aux performances qui peut être résolu automatiquement. Utilisez des solutions de AWS surveillance telles qu'[Amazon CloudWatch](#) ou AWS X-Ray pour vous aider à mieux comprendre la cause première du problème.
- Définissez le processus d'automatisation : créez un processus step-by-step de correction qui peut être utilisé pour résoudre automatiquement le problème.
- Configuration de l'événement d'initiation : configurez l'événement pour lancer automatiquement le processus de correction. Par exemple, vous pouvez définir un déclencheur pour redémarrer automatiquement une instance lorsqu'elle atteint un certain seuil d'CPU utilisation.
- Automatisez la correction : utilisez les AWS services et les technologies pour automatiser le processus de correction. Par exemple, [AWS Systems Manager Automation](#) fournit une solution

sécurisée et évolutive d'automatisation du processus de résolution. Veillez à utiliser une logique d'auto-réparation pour annuler les modifications si elles ne permettent pas de résoudre le problème.

- Test du flux de travail : testez le processus de résolution automatisé dans un environnement de pré-production.
- Mise en œuvre du flux de travail : implémentez la correction automatique dans l'environnement de production.
- Élaboration d'un manuel : élaborer et documenter un manuel qui décrit les étapes du plan de remédiation, y compris les événements initiateurs, la logique de remédiation et les mesures prises. Veillez à former les parties prenantes pour les aider à répondre efficacement aux événements de résolution automatisée.
- Révision et affinage : évaluez régulièrement l'efficacité du flux de travail de correction automatisé. Ajustez les événements de lancement et la logique de résolution, si nécessaire.

Ressources

Documents connexes :

- [CloudWatchDocumentation](#)
- [AWS Partner Network Partenaires de surveillance, de journalisation et de performance](#)
- [Documentation X-Ray](#)
- [Utilisation des alarmes et des actions d'alarme dans CloudWatch](#)
- [Élaborez une pratique d'automatisation du cloud pour l'excellence opérationnelle : les meilleures pratiques de AWS Managed Services](#)
- [Automatisez le réglage des performances de votre Amazon Redshift grâce à l'optimisation automatique des tables](#)

Vidéos connexes :

- [AWS re:Invent 2023 - Stratégies de mise à l'échelle automatisée, de correction et d'auto-réparation intelligente](#)
- [AWS re:Invent 2023 - \[LAUNCH\] Surveillance des applications pour les charges de travail modernes](#)
- [AWS re:Invent 2023 - Mise en œuvre de l'observabilité des applications](#)

- [AWS re:Invent 2021 - Automatisation intelligente des opérations dans le cloud](#)
- [AWS re:Invent 2022 - Configuration de contrôles à grande échelle dans votre environnement AWS](#)
- [AWS re:Invent 2022 - Automatisation de la gestion des correctifs et de la conformité à l'aide de AWS](#)
- [AWS re:Invent 2022 - Comment Amazon utilise de meilleurs indicateurs pour améliorer les performances de son site Web](#)
- [AWS re:Invent 2023 - Prenez le dessus : diagnostiquez et résolvez les problèmes de performance avec Amazon RDS](#)
- [AWS re:Invent 2021 - {New Launch} Détectez et résolvez automatiquement les problèmes avec Amazon Guru DevOps](#)
- [AWS re:Invent 2023 - Centralisez vos opérations](#)

Exemples connexes :

- [CloudWatch Journaux, personnalisation des alarmes](#)

PERF05-BP06 Maintenez votre charge de travail et vos services up-to-date

Restez up-to-date sur les nouveaux services et fonctionnalités du cloud pour adopter des fonctionnalités efficaces, résoudre les problèmes et améliorer l'efficacité globale des performances de votre charge de travail.

Anti-modèles courants :

- Vous supposez que votre architecture actuelle est statique et ne sera pas mise à jour au fil du temps.
- Vous ne disposez pas de systèmes ou de rythme régulier pour évaluer la compatibilité des packages et des logiciels mis à jour avec votre charge de travail.

Avantages de la mise en place de cette meilleure pratique : en établissant un processus pour rester à up-to-date jour avec les nouveaux services et offres, vous pouvez adopter de nouvelles fonctionnalités, résoudre les problèmes et améliorer les performances de la charge de travail.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : faible

Directives d'implémentation

Évaluez les méthodes d'amélioration des performances au fur et à mesure que de nouveaux services, modèles de conception et fonctionnalités de produits entrent en scène. Identifiez celles de ces méthodes qui sont susceptibles d'améliorer les performances ou d'accroître l'efficacité de la charge de travail via l'évaluation, la discussion interne ou l'analyse externe. Mettez en place un processus permettant d'évaluer les mises à jour, les nouvelles fonctions et les services pertinents pour votre charge de travail. Par exemple, créez une démonstration de faisabilité qui utilise les nouvelles technologies ou consultez un groupe interne. Lorsque vous essayez de nouvelles idées ou services, exécutez des tests de performances pour mesurer leur impact sur les performances de la charge de travail.

Étapes d'implémentation

- Inventaire de votre charge de travail : établissez l'inventaire de votre logiciel de charge de travail et de l'architecture, et identifiez les composants pouvant être mis à jour.
- Identification des sources de mise à jour : identifiez les actualités et mettez à jour les sources liées aux composants de votre charge de travail. Par exemple, vous pouvez vous abonner au [AWS blog What's New at](#) pour découvrir les produits correspondant à votre composante de charge de travail. Vous pouvez vous abonner au RSS flux ou gérer vos [abonnements par e-mail](#).
- Définition d'un calendrier de mise à jour : définissez un calendrier pour évaluer les nouveaux services et les nouvelles fonctionnalités adaptés à votre charge de travail.
 - Vous pouvez utiliser [AWS Systems Manager Inventory](#) pour collecter les métadonnées du système d'exploitation (OS), des applications et des instances à partir de vos EC2 instances Amazon et comprendre rapidement quelles instances exécutent le logiciel et les configurations requises par votre politique logicielle et quelles instances doivent être mises à jour.
- Évaluation de la nouvelle mise à jour : comprenez comment mettre à jour les composants de votre charge de travail. Profitez de l'agilité du cloud pour tester rapidement la façon dont les nouvelles fonctionnalités peuvent améliorer votre charge de travail afin de gagner en efficacité.
- Utiliser l'automatisation : utilisez l'automatisation pour le processus de mise à jour afin de réduire le niveau d'effort nécessaire au déploiement des nouvelles fonctionnalités et de limiter les erreurs causées par les processus manuels.
 - Vous pouvez utiliser [CI/CD](#) pour mettre à jour AMIs automatiquement des images de conteneur et d'autres artefacts liés à votre application cloud.

- Vous pouvez utiliser des outils tels que [AWS Systems Manager Patch Manager](#) pour automatiser le processus de mise à jour du système et planifier l'activité à l'aide de [AWS Systems Manager Maintenance Windows](#).
- Documentation du processus : documentez votre processus d'évaluation des mises à jour et des nouveaux services. Donnez aux propriétaires le temps et l'espace nécessaires pour rechercher, tester, expérimenter et valider les mises à jour et les nouveaux services. Reportez-vous aux exigences commerciales documentées et aidez KPIs à hiérarchiser les mises à jour qui auront un impact commercial positif.

Ressources

Documents connexes :

- [Blog AWS](#)
- [Quoi de neuf avec AWS](#)
- [up-to-dateImplémentation d'images avec des pipelines EC2 Image Builder automatisés](#)

Vidéos connexes :

- [AWS Re:inForce 2022 - Automatisation de la gestion des correctifs et de la conformité à l'aide de AWS](#)
- [All Things Patch : AWS Systems Manager | AWS Événements](#)

Exemples connexes :

- [Gestion de l'inventaire et des correctifs](#)
- [Un atelier sur l'observabilité](#)

PERF05-BP07 Vérification des métriques à intervalles réguliers

Vérifiez les métriques qui sont collectées dans le cadre de la maintenance de routine ou en réponse à des événements ou des incidents. Utilisez ces vérifications pour identifier d'une part les métriques qui ont été essentielles pour traiter les problèmes, et d'autre part les métriques supplémentaires, si elles ont été suivies, qui pourraient aider à identifier, traiter ou empêcher les problèmes.

Anti-modèles courants :

- Vous autorisez les métriques à rester dans un état d'alarme pendant longtemps.
- Vous créez des alarmes qui ne sont pas exploitables par un système d'automatisation.

Avantages liés au respect de cette bonne pratique : passez en revue en permanence les métriques qui sont collectées pour vérifier qu'elles identifient, résolvent ou préviennent correctement les problèmes. Les métriques peuvent également devenir caduques si vous les laissez dans un état d'alarme pendant longtemps.

Niveau d'exposition au risque si cette bonne pratique n'est pas respectée : moyen

Directives d'implémentation

Améliorez constamment la surveillance et la collecte des métriques. Lorsque vous répondez aux incidents ou aux événements, évaluez les métriques qui ont été utiles dans la gestion du problème et les métriques qui auraient pu aider mais ne sont pas suivies actuellement. Utilisez cette méthode pour améliorer la qualité des métriques que vous collectez afin de pouvoir prévenir ou résoudre plus rapidement les incidents futurs.

Lorsque vous répondez aux incidents ou aux événements, évaluez les métriques qui ont été utiles dans la gestion du problème et les métriques qui auraient pu aider mais ne sont pas suivies actuellement. Utilisez ce processus pour améliorer la qualité des métriques que vous collectez afin de pouvoir prévenir ou résoudre plus rapidement les incidents futurs.

Étapes d'implémentation

- **Définition de métriques** : définissez des métriques de performance critiques à surveiller qui correspondent à votre objectif de charge de travail, notamment des métriques telles que le temps de réponse et l'utilisation des ressources.
- **Établissement de bases de référence** : définissez une base de référence et une valeur souhaitable pour chaque métrique. La base de référence doit fournir des points de référence pour identifier les écarts ou les anomalies.
- **Établissement d'une cadence** : définissez une cadence (hebdomadaire ou mensuelle, par exemple) pour examiner les métriques critiques.
- **Identification des problèmes de performance** : au cours de chaque examen, évaluez les tendances et les écarts par rapport aux valeurs de référence. Recherchez les goulots d'étranglement ou les anomalies au niveau des performances. Pour les problèmes identifiés, effectuez une analyse détaillée des causes profondes afin de comprendre la raison principale du problème.

- Identification des actions correctives : utilisez votre analyse pour identifier les actions correctives. Cela peut inclure le réglage des paramètres, la correction de bogues et la mise à l'échelle des ressources.
- Documentation des résultats : documentez vos conclusions, y compris les problèmes identifiés, les causes profondes et les mesures correctives.
- Répétition et amélioration : évaluez et améliorez en permanence le processus de révision des métriques. Utilisez les enseignements tirés de la révision précédente pour améliorer le processus au fil du temps.

Ressources

Documents connexes :

- [Documentation CloudWatch](#)
- [Collecte de métriques et de journaux à partir d'instances Amazon EC2 et de serveurs sur site avec l'agent CloudWatch](#)
- [Interrogation de vos métriques avec CloudWatch Metrics Insights](#)
- [Surveillance, journalisation et performances : partenaires AWS Partner Network](#)
- [Documentation X-Ray](#)

Vidéos connexes :

- [AWS re:Invent 2022 - Setting up controls at scale in your AWS environment](#)
- [AWSre:Invent 2022 - How Amazon uses better metrics for improved website performance](#)
- [AWS re:Invent 2023 – Building an effective observability strategy](#)
- [AWS Summit SF 2022 - Full-stack observability and application monitoring with AWS](#)
- [AWSre:Invent 2023 - Take a load off: Diagnose & resolve performance issues with Amazon RDS](#)

Exemples connexes :

- [Création d'un tableau de bord avec QuickSight](#)
- [Tableaux de bord CloudWatch](#)

Conclusion

Pour atteindre et maintenir l'efficacité des performances, il est nécessaire d'avoir une approche axée sur les données. Vous devriez sérieusement envisager des modèles d'accès et des compromis qui vous permettront d'optimiser les performances. L'utilisation d'un processus d'évaluation basé sur des comparatifs et des tests de charge vous permet de sélectionner les configurations et types de ressources appropriés. En traitant votre infrastructure comme du code, vous pouvez faire évoluer votre architecture rapidement et en toute sécurité tout en utilisant les données pour prendre des décisions basées sur des faits en ce qui concerne votre architecture. La mise en place d'une surveillance à la fois active et passive permet de s'assurer que les performances de votre architecture ne se dégradent pas au fil du temps.

AWS s'efforce de vous aider à créer des architectures performantes tout en apportant une valeur commerciale. Utilisez les outils et techniques présentés dans ce document pour garantir votre réussite.

Collaborateurs

Les personnes et organisations suivantes ont contribué à l'élaboration du présent document :

- Sam Mokhtari, architecte principal de solutions en matière d'efficacité, Amazon Web Services
- Josh Hart, architecte de solutions, Amazon Web Services
- Richard Trabing, architecte de solutions, Amazon Web Services
- Brett Looney, architecte principal de solutions, Amazon Web Services
- Nina Vogl, architecte principal de solutions, Amazon Web Services
- Eric Pullen, architecte de solutions, Amazon Web Services
- Julien Lépine, responsable des architectes de solutions spécialisés, Amazon Web Services
- Ronnen Slasky, architecte de solutions, Amazon Web Services

Suggestions de lecture

Pour obtenir de l'aide, consultez les ressources suivantes :

- [Framework AWS Well-Architected](#)
- [Centre d'architecture AWS](#)

Révisions du document

Pour être informé des mises à jour de ce livre blanc, abonnez-vous au flux RSS.

Modification	Description	Date
Mise à jour mineure des bonnes pratiques	La bonne pratique PERF03-BP04 a été mise à jour avec de nouvelles recommandations de service.	6 novembre 2024
Mises à jour des conseils sur les bonnes pratiques	Plusieurs petites mises à jour dans l'ensemble du pilier.	27 juin 2024
Mise à jour et restructuration majeures	<p>Ce pilier a été restructuré pour inclure cinq domaines de bonnes pratiques (contre huit auparavant). Le contenu a été regroupé dans ces cinq domaines et a été mis à jour.</p> <p>Les nouveaux domaines de bonnes pratiques sont la sélection de l'architecture, le calcul et le matériel, la gestion des données, la mise en réseau et la diffusion de contenu, ainsi que les processus et la culture.</p>	3 octobre 2023
Mise à jour mineure	Suppression du langage non inclusif.	13 avril 2023
Mises à jour du nouveau cadre	Les bonnes pratiques ont été mises à jour avec des recommandations et de nouvelles bonnes pratiques.	10 avril 2023

Livre blanc mis à jour	Les bonnes pratiques ont été mises à jour avec de nouvelles recommandations en matière d'implémentation.	15 décembre 2022
Livre blanc mis à jour	Développement des bonnes pratiques et ajout de plans d'amélioration.	20 octobre 2022
Mise à jour mineure	Suppression du langage non inclusif.	22 avril 2022
Mises à jour mineures	Mise à jour des liens.	10 mars 2021
Mises à jour mineures	Le délai d'attente AWS Lambda a été modifié à 900 secondes et le nom d'Amazon Keyspaces (pour Apache Cassandra) a été corrigé.	5 octobre 2020
Mise à jour mineure	Correction d'un lien rompu.	15 juillet 2020
Mises à jour du nouveau cadre	Révision et mise à jour majeures du contenu	8 juillet 2020
Livre blanc mis à jour	Mise à jour mineure pour les problèmes grammaticaux	1er juillet 2018
Livre blanc mis à jour	Actualisation du livre blanc pour refléter les modifications apportées à AWS	1er novembre 2017
Publication initiale	Pilier Efficacité des performances – AWS Well-Architected Framework publié.	1er novembre 2016

Avis

Il incombe aux clients de procéder à une évaluation indépendante des informations contenues dans le présent document. Ce document : (a) est fourni à titre informatif uniquement, (b) représente les offres de AWS produits et les pratiques actuelles, qui sont susceptibles d'être modifiées sans préavis, et (c) ne crée aucun engagement ni aucune assurance de la part de AWS ses filiales, fournisseurs ou concédants de licence. AWS les produits ou services sont fournis « tels quels » sans garanties, déclarations ou conditions d'aucune sorte, qu'elles soient explicites ou implicites. Les responsabilités et obligations AWS de ses clients sont régies par AWS des accords, et ce document ne fait partie d'aucun accord conclu entre AWS et ses clients et ne les modifie pas.

© 2023, Amazon Web Services, Inc. ou ses sociétés apparentées. Tous droits réservés.

AWS Glossaire

Pour la AWS terminologie la plus récente, consultez le [AWS glossaire](#) dans la Glossaire AWS référence.