



Choix d'une base de données AWS vectorielle pour les cas d'utilisation de RAG

# AWS Conseils prescriptifs



---

# AWS Conseils prescriptifs: Choix d'une base de données AWS vectorielle pour les cas d'utilisation de RAG

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

Les marques et la présentation commerciale d'Amazon ne peuvent être utilisées en relation avec un produit ou un service qui n'est pas d'Amazon, d'une manière susceptible de créer une confusion parmi 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

Introduction .....	1
Public visé .....	1
Vue d'ensemble des vecteurs .....	3
Vue d'ensemble des bases de données vectorielles .....	5
Options de base de données vectorielles .....	7
Options de base de données vectorielles individuelles .....	7
Amazon Kendra .....	7
Amazon OpenSearch Service .....	8
Amazon RDS pour PostgreSQL avec pgvector .....	9
Amazon DocumentDB .....	9
Amazon MemoryDB .....	10
Amazon Neptune Analytics .....	11
Amazon S3 Vectors .....	12
Option de service géré .....	13
Choisir la bonne base de données vectorielle .....	14
Comparaison de bases de données vectorielles .....	16
Bases de données vectorielles individuelles .....	16
Service géré — Bases de connaissances Amazon Bedrock .....	20
Choisir entre des options individuelles et des options gérées .....	23
Comparaisons de coûts et considérations .....	25
Cas d'utilisation de bases de données vectorielles .....	30
Gestion des connaissances avec Amazon Kendra .....	30
Analyses en temps réel avec OpenSearch Serverless .....	31
Prochaines étapes et ressources .....	32
Ressources .....	32
AWS articles de blog .....	33
AWS documentation de service .....	33
Autres AWS ressources .....	33
Autres ressources .....	33
Historique du document .....	35
Glossaire .....	36
# .....	36
A .....	37
B .....	40

---

C .....	42
D .....	46
E .....	50
F .....	53
G .....	55
H .....	56
I .....	58
L .....	60
M .....	61
O .....	66
P .....	69
Q .....	72
R .....	72
S .....	75
T .....	80
U .....	81
V .....	82
W .....	82
Z .....	83
.....	lxxxv

# Choix d'une base de données AWS vectorielle pour les cas d'utilisation de RAG

Mayuri Shinde, Ivan Cui et Anand Bukkapatnam Tirumala, Amazon Web Services

Mars 2026 ([historique du document](#))

Les bases de données vectorielles sont de plus en plus importantes pour les organisations qui mettent en œuvre des applications d'IA générative. Ces bases de données stockent et gèrent des vecteurs, qui sont des représentations numériques de données permettant de traiter du texte, des images et d'autres contenus de manière à saisir leur signification et leurs relations.

À mesure que les entreprises explorent les options des bases de données vectorielles AWS, elles doivent comprendre les fonctionnalités, les compromis et les meilleures pratiques des différentes solutions. Ce guide vous aide à comparer les magasins de vecteurs couramment utilisés AWS et à prendre des décisions éclairées quant aux options les mieux adaptées à vos besoins ou à votre [cas d'utilisation](#) spécifiques. Que vous mettiez en œuvre la génération augmentée de récupération (RAG), que vous développiez des systèmes de recommandation ou que vous développiez d'autres applications d'IA, ce guide fournit un cadre qui vous aidera à évaluer et à choisir une solution de base de données vectorielle.

## Public visé

Ce guide est destiné aux personnes occupant les rôles suivants :

- Scientifiques des données et ingénieurs en apprentissage automatique (ML) qui utilisent des bases de données vectorielles pour stocker et récupérer des données de grande dimension pour les modèles de machine learning.
- Ingénieurs de données qui conçoivent et mettent en œuvre des pipelines de données comprenant des bases de données vectorielles pour le stockage et le traitement de données de grande dimension.
- MLOps ingénieurs qui utilisent des bases de données vectorielles dans le cadre du pipeline ML pour stocker et diffuser les sorties de modèles ou les représentations intermédiaires.
- Ingénieurs logiciels qui intègrent des bases de données vectorielles dans des applications nécessitant des systèmes de recherche ou de recommandation de similarité.

- DevOps ingénieurs chargés du déploiement et de la maintenance des bases de données vectorielles dans les environnements de production, afin de garantir l'évolutivité et la fiabilité.
- Les chercheurs en IA qui utilisent des bases de données vectorielles pour stocker et analyser de grands ensembles de données d'intégrations ou de vecteurs de caractéristiques.
- Chefs de produits d'IA qui ont besoin de comprendre les capacités et les limites des bases de données vectorielles pour prendre des décisions éclairées concernant les fonctionnalités et l'architecture des produits.

# Vue d'ensemble des vecteurs

Les vecteurs sont des représentations numériques qui aident les machines à comprendre et à traiter les données. Dans le domaine de l'IA générative, ils répondent à deux objectifs principaux :

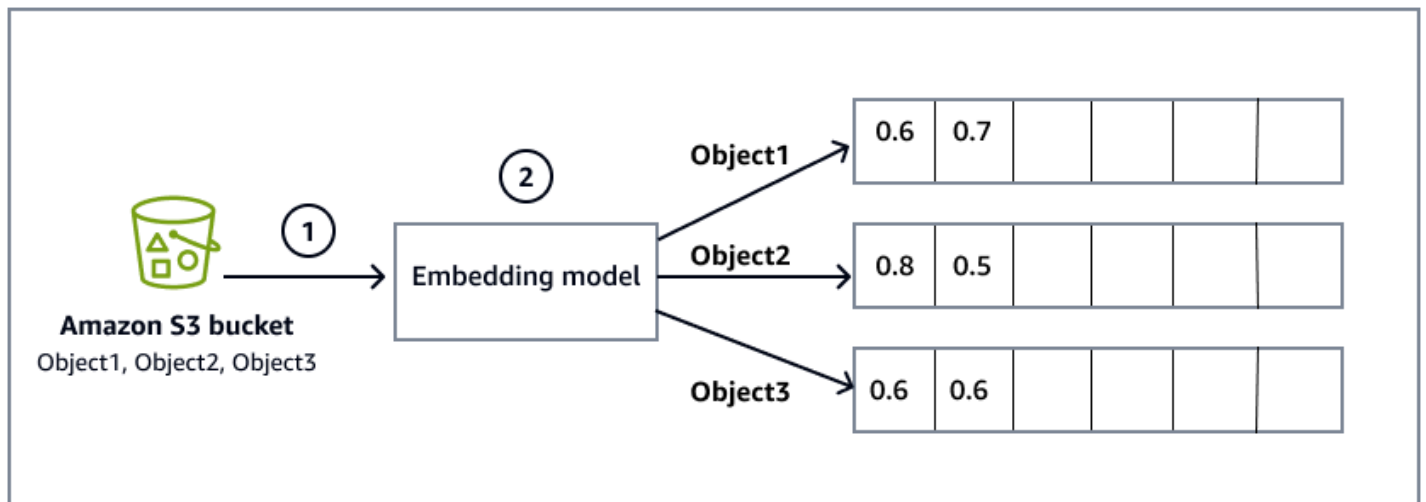
- Représentation des espaces latents qui capturent la structure des données sous forme compressée
- Création d'intégrations pour des données, telles que des mots, des phrases et des images

Les modèles d'intégration tels que [Word2Vec GloVe](#) et [Amazon Titan Text Embeddings](#) convertissent les données en vecteurs grâce à un processus appelé intégration. Ces modèles d'intégration peuvent effectuer les opérations suivantes :

- Apprenez à partir du contexte pour représenter les mots sous forme de vecteurs
- Rapprochez les mots similaires dans l'espace vectoriel
- Permettre aux machines de traiter les données dans un espace continu

Le schéma suivant fournit une vue d'ensemble détaillée du processus d'intégration :

1. Un bucket [Amazon Simple Storage Service \(Amazon S3\)](#) contient des fichiers qui sont les sources de données à partir desquelles le système lit et traite les informations. Le compartiment Amazon S3 est spécifié lors de la configuration de la base de connaissances [Amazon Bedrock](#), qui inclut également la [synchronisation des données avec la base de connaissances](#).
2. Le modèle d'intégration convertit les données brutes des fichiers objets du compartiment Amazon S3 en intégrations vectorielles. Par exemple, `Object1` est converti en un vecteur `[0.6, 0.7, ...]` qui représente son contenu dans un espace multidimensionnel.



Les intégrations de mots sont cruciales pour le traitement du langage naturel (NLP) car elles permettent :

- Capturez les relations sémantiques entre les mots
- Permettre la génération de texte contextuellement pertinent
- Alimenter de grands modèles linguistiques (LLMs) pour produire des réponses semblables à celles des humains

# Vue d'ensemble des bases de données vectorielles

Une base de données vectorielle est un système spécialisé qui stocke et interroge efficacement des vecteurs de grande dimension. Ces bases de données sont fondamentales pour les applications RAG (Retrieval Augmented Generation).

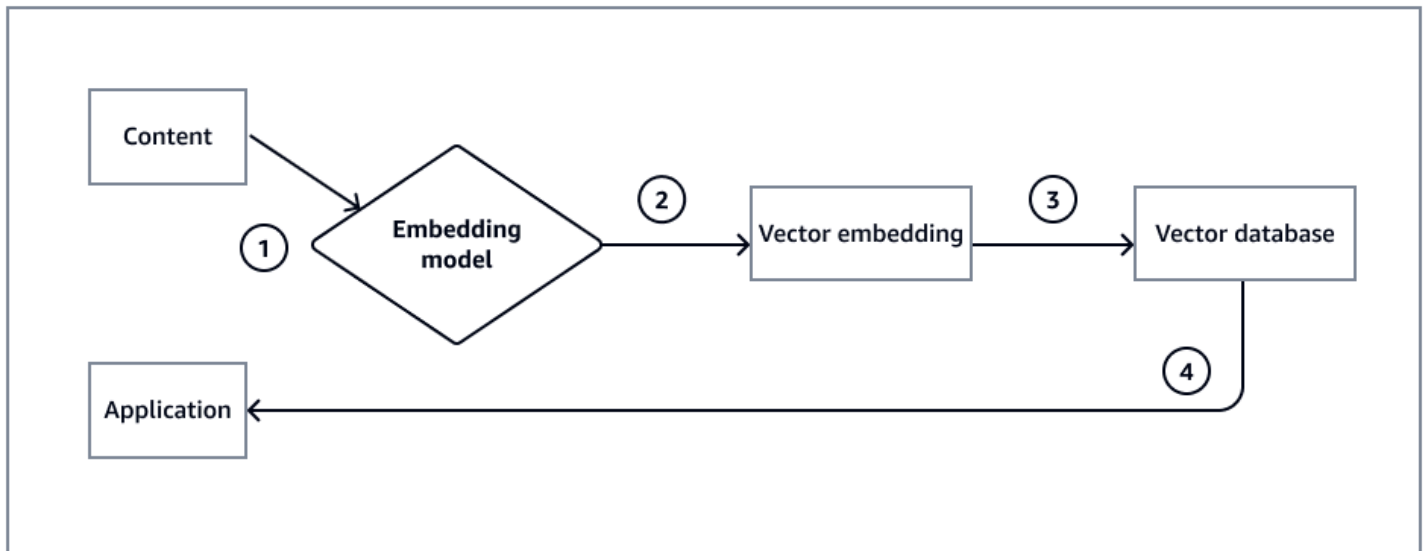
Les bases de données vectorielles gèrent la conversion et le stockage des données de la manière suivante :

- Les objets (tels que les fichiers audio, les images et les fichiers texte) sont convertis en vecteurs à l'aide de modèles d'intégration.
- Les vecteurs sont stockés dans des formats de données spécialisés.
- Les bases de données vectorielles permettent des recherches rapides de similarité.

Les bases de données vectorielles offrent plusieurs avantages essentiels par rapport aux bases de données traditionnelles, ce qui les rend particulièrement adaptées aux défis actuels en matière de données. Ils sont spécifiquement optimisés pour les opérations vectorielles et gèrent efficacement les données de grande dimension. Ils se spécialisent également dans les recherches de similarité auxquelles les bases de données traditionnelles ont du mal à faire face. Au-delà de ces fonctionnalités de base, les bases de données vectorielles sont conçues pour répondre aux demandes évolutives des applications de machine learning et d'IA générative. Ils excellent dans le stockage vectoriel à grande échelle et utilisent l'informatique distribuée pour équilibrer les charges de travail sur plusieurs nœuds. Cela garantit l'évolutivité et les performances à mesure que les volumes de données augmentent.

Le schéma suivant montre une implémentation de RAG :

1. Le contenu, tel que les documents ou les fichiers texte, est introduit dans le modèle d'intégration sous forme de données brutes à traiter. PDFs
2. Le modèle d'intégration transforme les données brutes en vecteurs numériques, qui représentent la signification sémantique du contenu.
3. Les intégrations vectorielles générées sont stockées dans une base de données vectorielle optimisée pour le stockage et la récupération de vecteurs de grande dimension.
4. Les applications peuvent désormais interroger la base de données vectorielle en réponse à des cas d'utilisation tels que la recherche sémantique et la recommandation de contenu.



Le choix d'une base de données vectorielle inappropriée pour une solution RAG peut entraîner des difficultés et des limites importantes, notamment les suivantes :

- Mauvaises performances des requêtes
- Les goulets d'étranglement liés à l'évolutivité
- Les défis de l'ingestion de données
- Absence de fonctionnalités avancées, telles que le filtrage et le classement
- Difficultés d'intégration avec d'autres systèmes
- Problèmes de persistance et de durabilité
- Problèmes de simultanéité et de cohérence dans les environnements comportant plusieurs utilisateurs
- Coûts de licence plus élevés ou dépendance vis-à-vis d'un fournisseur
- Soutien et ressources communautaires limités
- Risques potentiels en matière de sécurité et de conformité

# Options de base de données vectorielles

AWS propose une gamme variée de solutions de bases de données vectorielles pour répondre à différents cas d'utilisation et exigences dans les applications d'IA générative. Ces options peuvent être classées globalement en services de base de données individuels et en offres de services gérés, chacune présentant des caractéristiques et des avantages distincts. Comprendre ces options est essentiel pour les entreprises qui cherchent à mettre en œuvre efficacement des fonctionnalités de recherche vectorielle tout en maintenant des performances, une évolutivité et une rentabilité optimales.

Pour plus d'informations sur les solutions de base de données vectorielles, consultez les sections suivantes :

- [Options de base de données vectorielles individuelles](#)
- [Option de service géré](#)
- [Choisir la bonne base de données vectorielle](#)

## Options de base de données vectorielles individuelles

[Les options de base de données vectorielles individuelles disponibles AWS incluent Amazon Kendra, Amazon OpenSearch Service, AmazonRDS for PostgreSQL avec pgvector, Amazon MemoryDB, Amazon DocumentDB, Amazon DocumentDB, Amazon Neptune Analytics et Amazon S3 Vector.](#)

(Extension open source, pgvector ajoute la possibilité de stocker et de rechercher des intégrations vectorielles générées par ML.) Ces solutions proposent différentes approches de la recherche vectorielle, permettant aux entreprises de choisir en fonction de leur infrastructure existante, de leurs exigences techniques et de leurs [cas d'utilisation](#) spécifiques.

### Amazon Kendra

Amazon Kendra est un service de recherche intelligent destiné aux entreprises qui utilise le traitement du langage naturel et des algorithmes d'apprentissage automatique avancés pour renvoyer des réponses spécifiques aux questions de recherche à partir de vos données. Amazon Kendra simplifie la mise en œuvre de la fonctionnalité de recherche, ce qui en fait une solution backend efficace pour les applications d'IA générative.

Les autres fonctionnalités clés d'Amazon Kendra sont les suivantes :

- Connexions natives à plus de [40 sources de données](#)
- Fonctionnalités intégrées de préparation des données
- Configuration rapide ne nécessitant pas d'expertise technique approfondie

Les avantages d'Amazon Kendra sont les suivants :

- Traitement automatisé des données (découpage, ingestion, extraction)
- De puissantes options de personnalisation :
  - [Recherche par facettes](#)
  - [Analyses de recherche](#)
  - [Optimisation de la pertinence des recherches](#)
- Accès programmatique simple via le [AWS SDK pour Python \(Boto3\)](#)

Pour plus d'informations, consultez la section [Avantages d'Amazon Kendra](#) dans la documentation Amazon Kendra.

## Amazon OpenSearch Service

Amazon OpenSearch Service est un service géré qui vous aide à déployer, exploiter et dimensionner des clusters OpenSearch de services dans le AWS Cloud.

Les principales fonctionnalités du OpenSearch Service sont les suivantes :

- Moteur de recherche et d'analyse open source
- Architecture distribuée
- Traitement des données en temps réel

Certains avantages de l'utilisation du OpenSearch Service sont les suivants :

- Scalabilité horizontale
- RESTful Support de l'API
- Gère les données structurées et non structurées
- Analyse des données en temps réel
- Adapté à différentes tailles de déploiement

Pour plus d'informations, consultez la section [Fonctionnalités d'Amazon OpenSearch Service](#) dans la documentation du OpenSearch service.

## Amazon RDS pour PostgreSQL avec pgvector

Amazon RDS for [PostgreSQL](#) avec pgvector AWS associe le service de base de données relationnelle gérée à l'extension de traitement vectoriel de PostgreSQL. Cette combinaison permet aux entreprises de stocker et d'interroger des vecteurs de grande dimension tout en gérant Amazon RDS. La solution est particulièrement adaptée aux applications d'intelligence artificielle génératives qui nécessitent des opérations vectorielles en temps réel sans la surcharge liée à la gestion de l'infrastructure de base de données.

Les principaux avantages d'Amazon RDS pour PostgreSQL avec pgvector sont les suivants :

- Haute disponibilité
- Basculement automatique
- Rentable (pay-per-use)
- Surveillance intégrée
- Intégration de données vectorielles en temps réel

Pour plus d'informations, consultez les [avantages d'Amazon RDS](#) dans la documentation Amazon RDS.

## Amazon DocumentDB

Amazon DocumentDB (compatible avec MongoDB) est une base de données de documents qui offre des fonctionnalités de recherche vectorielle natives dans les versions 5.0 et ultérieures. Il combine la flexibilité du stockage de documents basé sur JSON avec la recherche vectorielle, prenant en charge à la fois les méthodes d'indexation hiérarchiques navigables Small World (HNSW) et Inverted File Flat (). IVFFlat

Les principales fonctionnalités d'Amazon DocumentDB sont les suivantes :

- Stockez et indexez des vecteurs jusqu'à 2 000 dimensions (jusqu'à 16 000 dimensions sans indexation)
- Temps de réponse en millisecondes pour les recherches de similarité vectorielle

- Support pour les mesures de distance entre produits euclidiens, cosinus et points
- Intégration parfaite avec les applications compatibles MongoDB existantes

Utilisez Amazon DocumentDB dans les situations suivantes :

- Pour les applications qui utilisent déjà MongoDB APIs et qui ont besoin de fonctionnalités de recherche vectorielle
- Pour les cas d'utilisation nécessitant des structures de données documentaires flexibles associées à une recherche sémantique
- Pour les scénarios nécessitant à la fois des requêtes documentaires traditionnelles et des recherches de similarité vectorielle
- Pour les applications proposant des recommandations de produits, des personnalisations, des assistants de chat et des services de détection des fraudes

Pour plus d'informations, consultez la section [Recherche vectorielle pour Amazon DocumentDB](#) dans la documentation Amazon DocumentDB.

## Amazon MemoryDB

Amazon MemoryDB est une base de données en mémoire compatible Redis qui fournit les performances de recherche vectorielle les plus rapides parmi les bases de données vectorielles les plus populaires sur AWS. Elle fournit des latences de requête inférieures à la milliseconde avec une durabilité dans les zones de multidisponibilité.

Les principales fonctionnalités de MemoryDB sont les suivantes :

- Stockez les données des applications et des millions de vecteurs dans une seule base de données
- Temps de réponse aux requêtes et aux mises à jour à un chiffre en millisecondes
- Les taux de rappel les plus élevés pour les performances les plus rapides sur AWS
- Support pour un maximum de 32 768 dimensions par vecteur
- Fonctionnalités de recherche sémantique et de mise en cache en temps réel

Utilisez MemoryDB dans les situations suivantes :

- Pour les applications en temps réel qui nécessitent une latence très faible (inférieure à 10 ms)

- Pour les charges de travail à haut débit avec des millions de demandes par jour
- Pour les cas d'utilisation tels que les moteurs de recommandation en temps réel, la mise en cache sémantique et la détection d'anomalies
- Pour les applications nécessitant à la fois un stockage de données en mémoire et des fonctionnalités de recherche vectorielle

Pour plus d'informations, consultez la section [Recherche vectorielle](#) dans la documentation de MemoryDB.

## Amazon Neptune Analytics

Amazon Neptune Analytics est un moteur d'analyse graphique qui offre des fonctionnalités de recherche vectorielle natives, ce qui le rend idéal pour les cas d'utilisation de Graph Retrieval Augmented Generation (GraphRag). Il combine la recherche de similarité vectorielle avec des traversées de graphes et des algorithmes.

Les principales fonctionnalités de Neptune Analytics sont les suivantes :

- Analysez des dizaines de milliards de relations en quelques secondes
- Combinez la recherche vectorielle avec des algorithmes graphiques (recherche de trajectoire, détection de communauté, centralité)
- Support pour les applications GraphRag avec des connaissances topologiques
- Jusqu'à 80 fois plus rapide que les solutions d'analyse graphique existantes
- Intégration à Amazon Bedrock pour une gestion complète de GraphRag

Utilisez Neptune Analytics dans les situations suivantes :

- Pour les applications GraphRag qui nécessitent des graphes de connaissances avec intégrations vectorielles
- Pour les cas d'utilisation qui nécessitent de parcourir des relations complexes parallèlement à la similarité vectorielle
- Pour les applications qui nécessitent des réponses explicables de l'IA avec un contexte relationnel
- Pour des scénarios tels que les vues à 360 degrés des clients, les réseaux de détection des fraudes et la découverte de connaissances

Pour plus d'informations, consultez la [documentation Amazon Neptune Analytics](#).

## Amazon S3 Vectors

Amazon S3 Vectors est le premier magasin d'objets cloud doté AWS de capacités natives de stockage vectoriel et de requêtes. Il fournit un stockage vectoriel spécialement conçu et optimisé en termes de coûts pour les applications d'IA nécessitant une grande échelle.

Les principales fonctionnalités d'Amazon S3 Vectors sont les suivantes :

- Stockage pour jusqu'à 2 milliards de vecteurs par index avec prise en charge de jusqu'à 10 000 index par compartiment de vecteurs
- Latence de requête inférieure à 100 ms optimisée pour le stockage à long terme et les modèles d'accès peu fréquents
- Jusqu'à 90 % de réduction des coûts pour les opérations vectorielles par rapport aux bases de données vectorielles spécialisées
- Architecture sans serveur avec mise à l'échelle automatique et durabilité de 99,999999999 % (11 s)

Utilisez les vecteurs Amazon S3 dans les situations suivantes :

- Pour les applications qui nécessitent le stockage de milliards de vecteurs à un coût minimal
- Pour les charges de travail qui tolèrent une latence de requête inférieure à une seconde (100 ms ou plus) au lieu de 10 ms
- Pour la conservation des vecteurs à long terme et les cas d'utilisation de l'archivage
- Pour les applications RAG avec des modèles de récupération peu fréquents
- Pour les entreprises qui privilégient l'économie du stockage par rapport à une latence extrêmement faible

Amazon S3 Vectors s'intègre nativement aux bases de connaissances Amazon Bedrock et fonctionne bien dans les architectures à plusieurs niveaux avec Amazon Service. OpenSearch Vous pouvez utiliser les vecteurs Amazon S3 pour le stockage à froid et utiliser le OpenSearch service pour les requêtes chaudes.

Pour plus d'informations, consultez la section [Utilisation des vecteurs S3 et des compartiments vectoriels](#) dans la documentation Amazon S3.

## Option de service géré

Amazon Bedrock Knowledge Bases représente l'approche AWS entièrement gérée de la mise en œuvre de bases de données vectorielles. La flexibilité des options de stockage du service, combinée à ses fonctionnalités de gestion automatisée, le rend particulièrement utile pour les entreprises qui cherchent à mettre en œuvre le RAG sans gérer une infrastructure complexe.

Avec les bases de connaissances Amazon Bedrock, vous pouvez créer, gérer et consulter des bases de connaissances qui améliorent vos modèles de base à l'aide de RAG. Ce service simplifie le processus complexe de mise en œuvre de RAG en gérant l'intégralité du pipeline d'ingestion, de vectorisation et de récupération des données.

Les principaux avantages des bases de connaissances Amazon Bedrock sont les suivants :

- Traitement des données simplifié
  - Ingestion et segmentation automatiques des données
  - Extraction de texte intégrée à partir de plusieurs formats de fichiers
  - Génération d'intégrations vectorielles gérées
  - Extraction et indexation automatiques des métadonnées
- Implémentation rationalisée du RAG
  - Stratégies de récupération préconfigurées
  - Optimisation automatique des fenêtres contextuelles
  - Réglage de la pertinence intégré
  - Fonctionnalités de recherche sémantique prêtes à l'emploi
- Sécurité et gouvernance
  - Contrôles intégrés Gestion des identités et des accès AWS (IAM)
  - Chiffrement des données au repos et en transit
  - Prise en charge de VPC
  - Journalisation des audits avec AWS CloudTrail

Les bases de connaissances Amazon Bedrock prennent en charge plusieurs [options de boutiques vectorielles](#), notamment :

- ~~Amazon Aurora PostgreSQL avec pgvector~~

- Amazon Neptune Analytics
- Amazon EMR sans serveur
- Amazon S3 Vectors
- Pinecone
- Redis Enterprise Cloud

Ce service géré gère l'ingestion, la vectorisation et la récupération automatisées des données. Cela simplifie les implémentations de RAG.

Pour obtenir des informations détaillées sur chaque boutique vectorielle prise en charge, consultez la [documentation des bases de connaissances Amazon Bedrock](#).

## Choisir la bonne base de données vectorielle

Sélectionnez votre base de données vectorielle en fonction de ces facteurs de décision clés :

- Si vous avez besoin d'une base de données de documents compatible avec MongoDB avec recherche vectorielle, choisissez Amazon DocumentDB. C'est idéal lorsque votre application utilise MongoDB APIs et que vous souhaitez ajouter des fonctionnalités de recherche sémantique sans gérer une infrastructure vectorielle distincte.
- Si vous avez besoin d'une latence très faible pour les applications en temps réel, choisissez Amazon MemoryDB. Cela permet d'obtenir les performances de recherche vectorielle les plus rapides AWS avec des temps de réponse inférieurs à la milliseconde. Il est idéal pour les moteurs de recommandation en temps réel et les applications à haut débit.
- Si vous avez besoin de représentations de connaissances basées sur des graphes avec recherche vectorielle, choisissez Amazon Neptune Analytics. C'est la solution idéale pour les applications GraphRag qui doivent traverser des relations complexes et effectuer des requêtes basées sur des graphes parallèlement à des recherches vectorielles, fournissant ainsi des réponses explicables par l'IA.
- Si vous devez associer des requêtes relationnelles à une recherche vectorielle, choisissez Amazon Aurora PostgreSQL avec pgvector. Cette option est idéale lorsque votre application nécessite à la fois des opérations SQL traditionnelles et des recherches de similarité vectorielle au sein de la même base de données.

- 
- Si vous avez besoin de requêtes à haut débit avec une latence inférieure à 10 ms, choisissez Amazon Service. OpenSearch II excelle dans la gestion des requêtes à haute fréquence et des applications en temps réel et inclut de récentes améliorations de l'accélération du GPU.
  - Si vous devez stocker des milliards de vecteurs de manière rentable, optez pour Amazon S3 Vectors. Cette option permet de réaliser jusqu'à 90 % d'économies et est idéale pour les applications dont les modèles de récupération sont peu fréquents (de quelques minutes à quelques heures entre les requêtes) et qui peuvent tolérer une latence inférieure à 100 ms.
  - Si vous avez besoin d'une recherche en texte intégral associée à une recherche vectorielle, choisissez Amazon OpenSearch Service. Cette option combine de puissantes fonctionnalités de recherche en texte intégral et de recherche vectorielle sur une seule plateforme.

# Comparaison de bases de données vectorielles

AWS propose plusieurs approches pour implémenter des fonctionnalités de recherche vectorielle, allant des bases de données vectorielles individuelles aux bases de connaissances Amazon Bedrock, qui est un service entièrement géré. Lors de l'évaluation de ces options, les entreprises doivent prendre en compte divers aspects, notamment l'architecture, l'évolutivité, les capacités d'intégration, les caractéristiques de performance et les fonctionnalités de sécurité.

## Bases de données vectorielles individuelles

Le tableau suivant fournit un aperçu des principales fonctionnalités de plusieurs solutions de base de données vectorielles AWS individuelles, en mettant l'accent sur leurs architectures, leurs capacités de mise à l'échelle, leurs intégrations de sources de données et leurs caractéristiques de performance.

Fonctionnalité	Amazon Kendra	Amazon OpenSearch Service	Amazon RDS pour PostgreSQL avec pgvector	Amazon DocumentDB	Amazon MemoryDB	Amazon Neptune Analytics	Amazon S3 Vectors
Cas d'utilisation principal	Recherche d'entreprise et RAG	Recherche et analyse distribuées	Base de données relationnelle avec support vectoriel	Base de données de documents avec recherche vectorielle	Recherche vectoriel en mémoire en temps réel	Analyse graphique avec recherche vectorielle	Stockage vectoriel optimisé en termes de coûts
Architecture	Entièrement géré	Cluster distribué	Base de données relationnelle	Orienté vers les documents	Base de données en mémoire	Moteur d'analyse graphique	Stockage d'objets sans serveur

Modèle de données	Basé sur des documents	Documents JSON	Tables relationnelles	Documents JSON	Valeur-clé avec JSON	Graphe de propriétés	Stockage d'objets
Dimensions du vecteur	Géré automatiquement	Jusqu'à 16 000	Configurable	Jusqu'à 2 000 (indexé) ; 16 000 (non indexé)	Jusqu'à 32 768	Configurable	Jusqu'à 4 096
Méthodes d'indexation	Automatique	NSW, FIV	NSW, IVFFlat	NSW, IVFFlat	HNSW	Graphe et vecteurs natifs	Automatique
Métriques de distance	Automatique	Cosine, Euclidean, produit à points	Cosine, produit euclidien, produit intérieur	Cosine, Euclidean, produit à points	Cosine, produit euclidien, produit intérieur	Cosinus euclidien	Cosinus euclidien
Latence des requêtes	Sous-seconde	Moins de 10 ms (accéléré par le GPU)	10 à 100 ms	Milliseconde	Moins d'une milliseconde	Sous-seconde	Moins de 100 ms
Modèle d'échelle	Automatique	Horizontal (ajouter des nœuds)	Répliques verticales et répliques en lecture	Horizontal (ajouter des instances)	Vertical et répliques	Automatique	Automatique (sans serveur)

Vecteurs maximaux	Gérées	Des milliards (en fonction du cluster)	Millions (en fonction de l'instance)	Des millions par collection	Des millions par base de données	Milliards	2 milliards par indice ; 10 000 indices par compartiment
Débit	Élevée	Très élevé (milliers de QPS)	Medium	Élevée	Très élevé (millions de demandes par jour)	Élevée	Moyen (optimisé pour les requêtes peu fréquentes)
Durabilité des données	99,999999 999 % (11 9 s)	Configurable avec des répliques	99,99 % (Multi-AZ)	99,99 % (Multi-AZ)	99,99 % (Multi-AZ)	99,99 %	99,999999 999 % (11 9 s)
Modèle de cohérence	Éventuel	Éventuel (configurable)	Fort (ACIDE)	Éventuel	Solide	Solide	Solide
Fonctionnalités supplémentaires	40 connecteurs de données ou plus, NLP	Recherche en texte intégral, analyses, tableaux de bord	Requêtes SQL, transactions ACID	Compatibilité avec l'API MongoDB	Compatibilité avec l'API Redis, mise en cache	Algorithmes graphiques, traversées	Intégration avec Amazon S3, politiques de cycle de vie

Modèle de tarification	Paiement par requête et par stockage	Heures d'ouverture et stockage des instances	Heures d'ouverture et stockage des instances	Heures d'ouverture et stockage des instances	Heures d'ouverture et stockage des instances	Unités de capacité et de stockage	Stockage, requêtes et transfert de données
Optimisation des coûts	Basé sur l'utilisation	Instances réservées, auto-scaling	Instances réservées, Aurora Serverless	Instances réservées	Instances réservées	Scalabilité automatique	Jusqu'à 90 % d'économies par rapport à une solution spécialisée DBs
Idéal pour	Recherche d'entreprise avec configuration minimale	Requêtes à haut débit et à faible latence	Charges de travail SQL et vectorielles hybrides	Applications compatibles avec MongoDB nécessitant des vecteurs	Applications en temps réel à très faible latence	GraphRag et graphes de connaissances	Stockage rentable et à long terme
Modèle de requête idéal	Recherches fréquentes dans les entreprises	Requêtes en temps réel à haute fréquence	Requêtes SQL et vectorielles mixtes	Requêtes de documents avec recherche sémantique	Des millions de demandes par jour	Parcours de graphes avec recherche vectorielle	Demandes peu fréquentes (de quelques minutes à quelques heures)

Complexité de configuration	Faible (entièrement géré)	Moyen (configuration en cluster)	Moyen (configuration de l'extension)	Moyen (configuration en cluster)	Moyen (configuration en cluster)	Faible (entièrement géré)	Faible (sans serveur)
Expertise d'équipe requise	Minimale	OpenSearch ou Elasticsearch	PostgreSQL, SQL	MongoDB	Redis	Bases de données graphiques	Amazon S3, concepts vectoriels de base

## Service géré — Bases de connaissances Amazon Bedrock

Les bases de connaissances Amazon Bedrock fournissent une solution entièrement gérée avec plusieurs options de stockage vectoriel. Le tableau suivant compare ces options de stockage.

Fonctionnalité	Aurora PostgreSQL avec SQL avec	Neptune Analytics	OpenSearch Service sans serveur	Vecteurs Amazon S3	Pomme de pin	RedisEnteprise Cloud
Cas d'utilisation principal	Base de données relationnelle avec vecteur RAG	Recherche vectorielle basée sur des graphes pour GraphRag	Gestion des connaissances RAG	RAG vectoriel optimisé en termes de coûts	Recherche vectorielle performante	Recherche vectorielle en mémoire
Architecture	Relationnel entièrement géré	Analyses graphiques entièrement gérées	Entièrement géré sans serveur	Stockage d'objets sans serveur	Cloud hybride entièrement géré	Entièrement géré en mémoire

Modèle de données	Tables relationnelles	Graphe de propriétés	Documents JSON	Stockage d'objets	Vecteurs spécialement conçus	Valeur-clé avec vecteurs
Stockage vectoriel	Grâce à l'extension pgvector	Vecteurs graphiques natifs	À travers le OpenSearch moteur	Stockage vectoriel natif d'Amazon S3	Base de données vectorielle native	Vecteurs en mémoire
Intégration avec Amazon Bedrock	Natif	Natif	Natif	Natif	Natif	Natif
Ingestion automatique	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)
Vectorisation automatique	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)	Oui (via Amazon Bedrock)
Mise à l'échelle	Mise à l'échelle automatique (Aurora Serverless)	Mise à l'échelle automatique des graphiques	Sans serveur automatique	Automatique (milliards de vecteurs)	Capsules à dimensionnement automatique	Clusters à dimensionnement automatique
Les performances des requêtes	Haut pour le mode relationnel ou vectoriel	Haut pour les vecteurs de graphes	Élevée	Moyen (latence de 100 ms ou plus)	Très élevée	Très élevée

Vecteurs maximaux	Millions (en fonction de l'instance)	Milliards	Milliards	2 milliards par indice	Milliards	Millions (en fonction de la mémoire)
Fonctionnalités supplémentaires	Requêtes SQL, transactions ACID	Algorithmes graphiques, traversées	Recherche en texte intégral, analyse	Cycle de vie d'Amazon S3, hiérarchisation	Filtrage des métadonnées, espaces de noms	Structures de données Redis, mise en cache
Optimisation des coûts	Modéré (Aurora Serverless)	Modéré (unités de capacité)	Élevé (sans serveur, pay-per-use)	Très élevé (jusqu'à 90 % d'économies)	Modéré (tarification basée sur les doses)	Faible (mémoire haut de gamme)
Idéal pour	Charges de SQL/vecteur travail hybrides	Graphes de connaissances connectés	Texte intégral avec recherche vectorielle	Vecteurs à long terme et à accès peu fréquent	Recherche vectorielle en temps réel à grande échelle	Besoins de latence très faible
Modèle de requête idéal	Requêtes SQL et vectorielles mixtes	Parcours de graphes à l'aide de vecteurs	Recherches fréquentes avec outils d'analyse	Récupération peu fréquente (de quelques minutes à quelques heures)	Requêtes en temps réel à haute fréquence	Des millions de demandes par seconde

Configuration avec Amazon Bedrock	Simple (géré par Amazon Bedrock)	Simple (géré par Amazon Bedrock)	Simple (géré par Amazon Bedrock)	Simple (géré par Amazon Bedrock)	Simple (géré par Amazon Bedrock)	Simple (géré par Amazon Bedrock)
Résidence des données	Régions AWS	Régions AWS	Régions AWS	Régions AWS	Multi-cloud (AWS et autres)	Multi-cloud (AWS et autres)
Modèle de tarification	Heures d'ouverture et stockage des instances	Unités de capacité et de stockage	Calcul et stockage (sans serveur)	Stockage, requêtes et transfert	Heures d'ouverture et stockage du pod	Heures d'ouverture et stockage des nœuds

## Choisir entre des options individuelles et des options gérées

Considération	Choisissez une base de données vectorielle individuelle	Choisissez les bases de connaissances Amazon Bedrock (gérées)
Implémentation du RAG	Vous voulez un contrôle total sur le pipeline RAG	Vous voulez un RAG entièrement géré avec une configuration minimale
Personnalisation	Vous avez besoin d'une logique de récupération et d'un prétraitement personnalisés	Les modèles RAG standard répondent à vos besoins
Infrastructures existantes	La base de données est déjà déployée	Vous reprenez un nouveau départ ou souhaitez une gestion simplifiée
Expertise de l'équipe	Votre équipe possède une expertise en administration de bases de données	Vous préférez vous concentrer sur la logique des applications plutôt que sur l'infrastructure

---

Complexité d'intégration	Vous avez besoin d'une intégration approfondie avec les systèmes existants	Vous souhaitez une intégration rapide avec les modèles Amazon Bedrock
Frais généraux d'exploitation	Vous pouvez gérer les opérations de base de données	Vous souhaitez AWS gérer les opérations
Structure des coûts	Vous préférez la tarification directe des bases de données	Vous préférez une tarification Amazon Bedrock unifiée
Délai de mise sur le marché	Vous avez le temps de procéder à une mise en œuvre personnalisée	Vous avez besoin d'un déploiement rapide

## Comparaisons de coûts et considérations

Comprendre la structure des coûts des différentes options de base de données vectorielles est essentiel pour prendre des décisions de mise en œuvre éclairées. Le tableau suivant décrit certaines considérations financières clés pour diverses solutions de base de données vectorielles, y compris les bases de données individuelles et les services gérés. Chaque option comporte des facteurs de tarification distincts qui peuvent avoir un impact sur le coût total de possession (TCO), qu'il s'agisse des pay-as-you-go modèles, de l'infrastructure ou des coûts opérationnels.

Base de données vectorielle	Modèle de coûts	Considérations relatives aux coûts
Amazon Kendra	Payez au fur et à mesure, en fonction des demandes	Les coûts peuvent varier en fonction du nombre de requêtes et de la quantité de données indexées. Des frais supplémentaires peuvent s'appliquer pour le stockage et le transfert de données. Pour plus d'informations, consultez les <a href="#">tarifs d'Amazon Kendra</a> .
Amazon OpenSearch Service	Payez au fur et à mesure, en fonction des heures d'utilisation de l'instance et du stockage	Les coûts incluent les heures d'instance, le stockage (volumes Amazon EBS), le transfert de données et le UltraWarm stockage optionnel . Les instances réservées peuvent permettre de réaliser jusqu'à 30 % d'économies. Les instances accélérées par GPU offrent un meilleur rapport prix/performances pour les charges de travail vectoriel les. Pour plus d'informations,

consultez les [tarifs d'Amazon OpenSearch Service](#).

Open source OpenSearch	Open source, sans frais directs (vous n'avez pas à payer pour télécharger ou utiliser le logiciel, et il n'y a aucun coût de licence)	Les coûts incluent l'infrastructure (tels que les serveurs et le stockage) et les coûts opérationnels (tels que la maintenance et la surveillance). Organisations doivent prévoir un budget pour le personnel nécessaire à la gestion et à la maintenance de l'infrastructure.
Amazon RDS pour PostgreSQL avec pgvector	Payez au fur et à mesure, en fonction de votre consommation	Les coûts incluent les types d'instances de base de données, le stockage, le transfert de données et les sauvegardes. Des frais supplémentaires peuvent s'appliquer pour le transfert de données, les types d'instances et le stockage au-delà du niveau AWS gratuit. Pour plus d'informations, consultez <a href="#">Tarification d'Amazon RDS</a> .

---

Amazon DocumentDB	Payez au fur et à mesure, en fonction des heures d'utilisation de l'instance et du stockage	Les coûts incluent les heures d'instance, le stockage (Go par mois), les I/O demandes, le stockage de sauvegarde et le transfert de données. Les clusters élastiques permettent une mise à l'échelle dynamique. Instances réservées disponibles pour l'optimisation des coûts. Pour plus d'informations, consultez la tarification <a href="#">d'Amazon DocumentDB</a> .
Amazon MemoryDB	Payez au fur et à mesure, en fonction des heures d'ouverture des nœuds et du stockage des données	Les coûts incluent les heures de nœud (par type de nœud), le stockage des données (Go par heure), le stockage des instantanés et le transfert de données. Les nœuds réservés peuvent permettre de réaliser jusqu'à 55 % d'économies. Optimisé pour les charges de travail à haut débit et à faible latence. Pour plus d'informations, consultez la tarification <a href="#">d'Amazon MemoryDB</a> .

---

Amazon Neptune Analytics	Payez au fur et à mesure, en fonction des unités de capacité	Les coûts incluent les unités de capacité Neptune (NCUs), le stockage (Go par mois) et le transfert de données. Mise à l'échelle automatique en fonction de la charge de travail, sans engagement initial. Un minimum de 128 NCUs est requis. Pour plus d'informations, consultez la <a href="#">tarification d'Amazon Neptune</a> .
Amazon S3 Vectors	Payez au fur et à mesure, en fonction de l'espace de stockage et des demandes	Les coûts incluent le stockage (Go par mois), les requêtes PUT et GET, la gestion des index vectoriels et le transfert de données. Permet de réaliser jusqu'à 90 % d'économies par rapport aux bases de données vectorielles spécialisées. Les politiques de hiérarchisation intelligente et de cycle de vie Amazon S3 sont disponibles pour une optimisation supplémentaire. Pour plus d'informations, consultez <a href="#">Tarification Amazon S3</a> .

---

Bases de connaissances  
Amazon Bedrock

Payez au fur et à mesure, en fonction de votre consommation

Les coûts peuvent varier en fonction de l'utilisation de la base de connaissances et de services supplémentaires tels qu'Amazon OpenSearch Serverless. Des frais supplémentaires peuvent s'appliquer pour le stockage des données, le transfert de données et les fonctionnalités supplémentaires. Pour plus d'informations sur les tarifs, consultez la section [Tarification d'Amazon OpenSearch Service](#).

# Cas d'utilisation de bases de données vectori

Les exemples suivants montrent comment différentes options de base de données vectorielles peuvent être utilisées efficacement pour améliorer la gestion des connaissances, améliorer l'efficacité opérationnelle et obtenir de meilleurs résultats commerciaux. Ces cas d'utilisation illustrent les applications pratiques des solutions de base de données vectorielles évoquées plus haut dans ce guide et fournissent un aperçu de leurs performances et avantages réels.

## Gestion des connaissances avec Amazon Kendra

**Problème avec le client** — L'un des plus grands entrepreneurs généraux du Japon était confronté à une baisse de son personnel expérimenté. L'entreprise avait besoin d'un moyen de transférer efficacement les connaissances et les compétences du personnel expérimenté à la jeune génération. Ils avaient besoin d'une solution pour saisir et diffuser les connaissances complexes en ingénierie de la construction et les expériences passées.

**AWS solution** — Pour résoudre ce problème, le client s'est tourné vers Amazon Kendra, une solution d'intelligence artificielle capable de gérer rapidement et précisément sa base de connaissances interne et d'autoriser les requêtes en langage naturel. Grâce à Amazon Kendra, les employés peuvent désormais trouver les informations dont ils ont besoin beaucoup plus rapidement, ce qui améliore la productivité et facilite le transfert de connaissances entre le personnel expérimenté et le personnel plus jeune.

**Impact** — En mettant en œuvre un chatbot génératif basé sur l'IA alimenté par Amazon Kendra, l'entreprise a créé une plateforme de connaissances unifiée. Le chatbot permet aux employés d'accéder rapidement aux connaissances techniques et aux expériences passées en génie de la construction. Cette solution a considérablement amélioré l'efficacité du transfert de connaissances et des processus de prise de décision au sein de l'organisation, contribuant ainsi à garantir que la précieuse expertise est préservée et facilement accessible à tous les employés. Le coût de cette solution peut varier en fonction de votre utilisation et de votre configuration. Pour une estimation détaillée des coûts, consultez le [Calculateur de tarification AWS](#). Pour l'estimation des coûts des bases de données vectorielles, consultez la section [Comparaison des coûts et considérations](#) de ce guide ou consultez les [tarifs d'Amazon Kendra](#).

Pour plus d'informations sur les autres cas d'utilisation, consultez les clients [d'Amazon Kendra](#).

## Analyses en temps réel avec OpenSearch Serverless

**Problème client** — Un fournisseur de services financiers de premier plan a dû relever le défi de gérer un énorme écosystème de données. Elle a traité 300 millions d'autorisations et 90 milliards de transactions par an, accumulant environ 1,1 pétaoctet (Po) de données. Le système existant, qui dessert 300 000 utilisateurs ayant besoin d'accéder à plus de 6 000 rapports, avait besoin d'être modernisé pour assurer une cohérence globale et permettre une prise de décision en temps réel.

**AWS solution** — L'architecture de la solution a utilisé des modèles de base disponibles via Amazon Bedrock (notamment Anthropic, Sonnet 3, Sonnet 3.5 et Haiku) pour le traitement du langage naturel. Le client a choisi OpenSearch Serverless comme base de données vectorielle pour son évolutivité supérieure et sa capacité à gérer efficacement l'énorme volume de données. Cette architecture a permis le traitement fluide des requêtes complexes et la génération de rapports dynamiques.

**Impact** — La mise en œuvre a permis d'augmenter la productivité de 50 % en éliminant le besoin de générer manuellement plus de 100 tableaux de bord de business intelligence. Les utilisateurs peuvent désormais générer des rapports par le biais de requêtes en langage naturel avec des temps de réponse compris entre 20 et 40 secondes. Le coût de cette solution peut varier en fonction de votre utilisation et de votre configuration. Pour une estimation détaillée des coûts, consultez le [Calculateur de tarification AWS](#). Pour l'estimation des coûts des bases de données vectorielles, consultez la section [Comparaison des coûts et considérations](#) de ce guide ou consultez les [tarifs d'Amazon OpenSearch Service](#). Pour plus d'informations sur les autres cas d'utilisation par les clients, consultez [Amazon OpenSearch Serverless](#).

## Prochaines étapes et ressources

Après avoir examiné ce guide, considérez les actions suivantes pour passer de la compréhension à la mise en œuvre :

1. Évaluez vos besoins actuels :
  - Évaluez votre infrastructure de base de données et votre expertise existantes.
  - Documentez vos exigences spécifiques en matière de recherche vectorielle.
  - Définissez vos objectifs de performance, de mise à l'échelle et de coûts.
2. Choisissez l'une des options suivantes pour tester les options de base de données vectorielles :
  - Option 1 : Configurez une preuve de concept à l'aide de votre solution de base de données vectorielle préférée.
  - Option 2 : testez des exemples de jeux de données dans les bases de connaissances Amazon Bedrock. Essayez l'expérience de création rapide pour une base de connaissances Amazon Bedrock. Par exemple, consultez la section [Création rapide d'une base de connaissances Aurora PostgreSQL pour Amazon Bedrock](#) dans la documentation Aurora.
3. Passez en revue [les ressources](#) supplémentaires.
4. Obtenez l'aide d'un expert :
  - Contactez votre Compte AWS équipe ou vos architectes de AWS solutions pour obtenir des conseils de mise en œuvre.
  - [Collaborez avec AWS des partenaires](#) spécialisés dans les bases de données vectorielles.
5. Planifiez votre déploiement en production :
  - Créez une stratégie de migration en cas de migration à partir de bases de données existantes.
  - Développez un plan de mise à l'échelle pour la solution que vous avez choisie.
  - Concevez vos procédures de surveillance et de maintenance.

## Ressources

Les ressources suivantes peuvent vous aider à choisir une base de données vectorielle.

## AWS articles de blog

- [Accélérez le développement de vos applications d'IA générative avec Amazon Bedrock Knowledge Bases, Quick Create et Amazon Aurora Serverless](#)
- [Explication des fonctionnalités de la base de données vectorielle d'Amazon OpenSearch Service](#)
- [Explorez en profondeur les magasins de données vectorielles à l'aide des bases de connaissances Amazon Bedrock](#)
- [Tirez parti de pgvector et Amazon Aurora PostgreSQL pour le traitement du langage naturel, les chatbots et l'analyse des sentiments](#)

## AWS documentation de service

- [Choix d'un service AWS de base de données](#)
- [Comment fonctionnent les bases de connaissances Amazon Bedrock](#)
- [Documentation de Neptune Analytics](#)
- [Présentation d'Amazon Web Services : bases de données](#)
- [Utilisation d'Aurora PostgreSQL comme base de connaissances pour Amazon Bedrock](#)
- [Utilisation d'Amazon Aurora PostgreSQL](#)
- [Amazon DocumentDB](#)
- [Amazon MemoryDB](#)
- [Amazon S3 Vectors](#)

## Autres AWS ressources

- [Bases de connaissances Amazon Bedrock](#)
- [Bases de données vectorielles et intégrations](#)
- [Bases de données vectorielles pour les applications d'IA générative](#)
- [Que sont les intégrations dans le Machine Learning ?](#)

## Autres ressources

- [À propos de PostgreSQL](#)

- 
- [documentation pgvector](#)
  - [Pinecone comme base de connaissances pour Amazon Bedrock](#)
  - [Redis Enterprise Cloud sur AWS](#)

## Historique du document

Le tableau suivant décrit les modifications importantes apportées à ce guide, intitulé Choisir une base de données AWS vectorielle pour les cas d'utilisation de RAG. Pour être averti des mises à jour à venir, abonnez-vous à un [fil RSS](#).

Modification	Description	Date
<a href="#">Ajouté Services AWS</a>	Nous avons ajouté des informations sur l'utilisation d'Amazon DocumentDB, d'Amazon MemoryDB, d'Amazon S3 Vectors et d'Amazon Neptune Analytics.	30 mars 2026
<a href="#">Publication initiale</a>	—	6 mars 2025

# AWS Glossaire des directives prescriptives

Les termes suivants sont couramment utilisés dans les stratégies, les guides et les modèles fournis par les directives AWS prescriptives. Pour suggérer des entrées, veuillez utiliser le lien [Faire un commentaire](#) à la fin du glossaire.

## Nombres

### 7 R

Sept politiques de migration courantes pour transférer des applications vers le cloud. Ces politiques s'appuient sur les 5 R identifiés par Gartner en 2011 et sont les suivantes :

- **Refactor/re-architect** — Déplacez une application et modifiez son architecture en tirant pleinement parti des fonctionnalités natives du cloud pour améliorer l'agilité, les performances et l'évolutivité. Cela implique généralement le transfert du système d'exploitation et de la base de données. Exemple : migrez votre base de données Oracle sur site vers l' PostgreSQL-Compatible édition Amazon Aurora.
- **Replatformer (déplacer et remodeler)** : transférez une application vers le cloud et introduisez un certain niveau d'optimisation pour tirer parti des fonctionnalités du cloud. Exemple : migrez votre base de données Oracle sur site vers Amazon Relational Database Service (Amazon RDS) pour Oracle dans le. AWS Cloud
- **Racheter (rachat)** : optez pour un autre produit, généralement en passant d'une licence traditionnelle à un modèle SaaS. Exemple : migrez votre système de gestion de la relation client (CRM) vers Salesforce.com.
- **Réhéberger (lift and shift)** : transférez une application vers le cloud sans apporter de modifications pour tirer parti des fonctionnalités du cloud. Exemple : migrez votre base de données Oracle sur site vers Oracle sur une instance EC2 dans le. AWS Cloud
- **Relocaliser (lift and shift au niveau de l'hyperviseur)** : transférez l'infrastructure vers le cloud sans acheter de nouveau matériel, réécrire des applications ou modifier vos opérations existantes. Vous migrez des serveurs d'une plateforme sur site vers un service cloud pour la même plateforme. Exemple : migrer une Microsoft Hyper-V application vers AWS.
- **Retenir** : conservez les applications dans votre environnement source. Il peut s'agir d'applications nécessitant une refactorisation majeure, que vous souhaitez retarder, et d'applications existantes que vous souhaitez retenir, car rien ne justifie leur migration sur le plan commercial.

- Retirer : mettez hors service ou supprimez les applications dont vous n'avez plus besoin dans votre environnement source.

## A

### A2 (1) Agent-to-Agent

Protocole dynamique pour la collaboration agent-agent prenant en charge la délégation de tâches et le transfert d'état.

### ABAC

Voir contrôle [d'accès basé sur les attributs](#).

### services abstraits

Consultez la section [Services gérés](#).

### ACIDE

Voir [atomicité, consistance, isolation, durabilité](#).

### migration active-active

Méthode de migration de base de données dans laquelle la synchronisation des bases de données source et cible est maintenue (à l'aide d'un outil de réplication bidirectionnelle ou d'opérations d'écriture double), tandis que les deux bases de données gèrent les transactions provenant de la connexion d'applications pendant la migration. Cette méthode prend en charge la migration par petits lots contrôlés au lieu d'exiger un basculement ponctuel. Elle est plus flexible mais demande plus de travail qu'une migration [active-passive](#).

### migration active-passive

Méthode de migration de base de données dans laquelle les bases de données source et cible sont synchronisées, mais seule la base de données source gère les transactions liées à la connexion des applications pendant que les données sont répliquées vers la base de données cible. La base de données cible n'accepte aucune transaction pendant la migration.

### Agent

Un système d'IA capable de raisonner, de planifier et de prendre des mesures de manière autonome à l'aide d'outils pour atteindre des objectifs.

## Agent Ops

Pratiques opérationnelles pour la création, le test, le déploiement et l'exécution d'agents d'IA en production à grande échelle.

### fonction d'agrégation

Fonction SQL qui agit sur un groupe de lignes et calcule une valeur de retour unique pour le groupe. Des exemples de fonctions d'agrégation incluent SUM et MAX.

## AI

Voir [intelligence artificielle](#).

### AIOps

Voir les [opérations d'intelligence artificielle](#).

### anonymisation

Processus de suppression définitive d'informations personnelles dans un ensemble de données. L'anonymisation peut contribuer à protéger la vie privée. Les données anonymisées ne sont plus considérées comme des données personnelles.

### anti-motif

Solution fréquemment utilisée pour un problème récurrent lorsque la solution est contre-productive, inefficace ou moins efficace qu'une solution alternative.

### contrôle des applications

Une approche de sécurité qui permet d'utiliser uniquement des applications approuvées afin de protéger un système contre les logiciels malveillants.

### portefeuille d'applications

Ensemble d'informations détaillées sur chaque application utilisée par une organisation, y compris le coût de génération et de maintenance de l'application, ainsi que sa valeur métier. Ces informations sont essentielles pour [le processus de découverte et d'analyse du portefeuille](#) et permettent d'identifier et de prioriser les applications à migrer, à moderniser et à optimiser.

### intelligence artificielle (IA)

Domaine de l'informatique consacré à l'utilisation des technologies de calcul pour exécuter des fonctions cognitives généralement associées aux humains, telles que l'apprentissage, la résolution de problèmes et la reconnaissance de modèles. Pour plus d'informations, veuillez consulter [Qu'est-ce que l'intelligence artificielle ?](#)

## opérations d'intelligence artificielle (AIOps)

Processus consistant à utiliser des techniques de machine learning pour résoudre les problèmes opérationnels, réduire les incidents opérationnels et les interventions humaines, mais aussi améliorer la qualité du service. Pour plus d'informations sur la façon dont les AIOps sont utilisées dans la stratégie de migration AWS, veuillez consulter le [guide d'intégration des opérations](#).

## chiffrement asymétrique

Algorithme de chiffrement qui utilise une paire de clés, une clé publique pour le chiffrement et une clé privée pour le déchiffrement. Vous pouvez partager la clé publique, car elle n'est pas utilisée pour le déchiffrement, mais l'accès à la clé privée doit être très restreint.

## atomicité, cohérence, isolement, durabilité (ACID)

Ensemble de propriétés logicielles garantissant la validité des données et la fiabilité opérationnelle d'une base de données, même en cas d'erreur, de panne de courant ou d'autres problèmes.

## contrôle d'accès par attributs (ABAC)

Pratique qui consiste à créer des autorisations détaillées en fonction des attributs de l'utilisateur, tels que le service, le poste et le nom de l'équipe. Pour plus d'informations, consultez [ABAC pour AWS](#) dans la documentation Gestion des identités et des accès AWS (IAM).

## source de données faisant autorité

Emplacement où vous stockez la version principale des données, considérée comme la source d'information la plus fiable. Vous pouvez copier les données de la source de données officielle vers d'autres emplacements à des fins de traitement ou de modification des données, par exemple en les anonymisant, en les expurgant ou en les pseudonymisant.

## Zone de disponibilité

Un emplacement distinct au sein d'une Région AWS réseau isolé des défaillances dans d'autres zones de disponibilité et fournissant une connectivité réseau peu coûteuse et à faible latence aux autres zones de disponibilité de la même région.

## AWS Cadre d'adoption du cloud (AWS CAF)

Un cadre de directives et de meilleures pratiques visant AWS à aider les entreprises à élaborer un plan efficace pour réussir leur migration vers le cloud. AWS La CAF organise ses conseils en six domaines prioritaires appelés perspectives : les affaires, les personnes, la gouvernance, les plateformes, la sécurité et les opérations. Les perspectives d'entreprise, de personnes et de gouvernance mettent l'accent sur les compétences et les processus métier, tandis que les

perspectives relatives à la plateforme, à la sécurité et aux opérations se concentrent sur les compétences et les processus techniques. Par exemple, la perspective liée aux personnes cible les parties prenantes qui s'occupent des ressources humaines (RH), des fonctions de dotation en personnel et de la gestion des personnes. Dans cette perspective, la AWS CAF fournit des conseils pour le développement du personnel, la formation et les communications afin de préparer l'organisation à une adoption réussie du cloud. Pour plus d'informations, veuillez consulter le [site Web AWS CAF](#) et le [livre blanc AWS CAF](#).

## AWS Cadre de qualification de la charge de travail (AWS WQF)

Outil qui évalue les charges de travail liées à la migration des bases de données, recommande des stratégies de migration et fournit des estimations de travail. AWS Le WQF est inclus avec AWS Schema Conversion Tool (AWS SCT). Il analyse les schémas de base de données et les objets de code, le code d'application, les dépendances et les caractéristiques de performance, et fournit des rapports d'évaluation.

## B

mauvais bot

Un [bot](#) destiné à perturber ou à nuire à des individus ou à des organisations.

BCP

Consultez la section [Planification de la continuité des activités](#).

graphique de comportement

Vue unifiée et interactive des comportements des ressources et des interactions au fil du temps. Vous pouvez utiliser un graphique de comportement avec Amazon Detective pour examiner les tentatives de connexion infructueuses, les appels d'API suspects et les actions similaires. Pour plus d'informations, veuillez consulter [Data in a behavior graph](#) dans la documentation Detective.

système de poids fort

Système qui stocke d'abord l'octet le plus significatif. Voir aussi [endianité](#).

classification binaire

Processus qui prédit un résultat binaire (l'une des deux classes possibles). Par exemple, votre modèle de machine learning peut avoir besoin de prévoir des problèmes tels que « Cet e-mail est-il du spam ou non ? » ou « Ce produit est-il un livre ou une voiture ? ».

## filtre de Bloom

Structure de données probabiliste et efficace en termes de mémoire qui est utilisée pour tester si un élément fait partie d'un ensemble.

## blue/green déploiement

Stratégie de déploiement dans laquelle vous créez deux environnements distincts mais identiques. Vous exécutez la version actuelle de l'application dans un environnement (bleu) et la nouvelle version de l'application dans l'autre environnement (vert). Cette stratégie vous permet de revenir rapidement en arrière avec un impact minimal.

## bot

Application logicielle qui exécute des tâches automatisées sur Internet et simule l'activité ou l'interaction humaine. Certains robots sont utiles ou bénéfiques, comme les robots d'exploration Web qui indexent des informations sur Internet. D'autres robots, appelés « bots malveillants », sont destinés à perturber ou à nuire à des individus ou à des organisations.

## botnet

Réseaux de [robots](#) infectés par des [logiciels malveillants](#) et contrôlés par une seule entité, connue sous le nom d'herder ou d'opérateur de bots. Les botnets sont le mécanisme le plus connu pour faire évoluer les bots et leur impact.

## branche

Zone contenue d'un référentiel de code. La première branche créée dans un référentiel est la branche principale. Vous pouvez créer une branche à partir d'une branche existante, puis développer des fonctionnalités ou corriger des bogues dans la nouvelle branche. Une branche que vous créez pour générer une fonctionnalité est communément appelée branche de fonctionnalités. Lorsque la fonctionnalité est prête à être publiée, vous fusionnez à nouveau la branche de fonctionnalités dans la branche principale. Pour plus d'informations, consultez [À propos des branches](#) (GitHub documentation).

## accès par brise-vitre

Dans des circonstances exceptionnelles et par le biais d'un processus approuvé, c'est un moyen rapide pour un utilisateur d'accéder à un accès auquel Compte AWS il n'est généralement pas autorisé. Pour plus d'informations, consultez l'indicateur [Mettre en œuvre des procédures permettant de briser le verre](#) dans le AWS Well-Architected guide.

## stratégie existante (brownfield)

L'infrastructure existante de votre environnement. Lorsque vous adoptez une stratégie existante pour une architecture système, vous concevez l'architecture en fonction des contraintes des systèmes et de l'infrastructure actuels. Si vous étendez l'infrastructure existante, vous pouvez combiner des politiques brownfield (existantes) et [greenfield](#) (inédites).

## cache de tampon

Zone de mémoire dans laquelle sont stockées les données les plus fréquemment consultées.

## capacité métier

Ce que fait une entreprise pour générer de la valeur (par exemple, les ventes, le service client ou le marketing). Les architectures de microservices et les décisions de développement peuvent être dictées par les capacités métier. Pour plus d'informations, veuillez consulter la section [Organisation en fonction des capacités métier](#) du livre blanc [Exécution de microservices conteneurisés sur AWS](#).

## planification de la continuité des activités (BCP)

Plan qui tient compte de l'impact potentiel d'un événement perturbateur, tel qu'une migration à grande échelle, sur les opérations, et qui permet à une entreprise de reprendre ses activités rapidement.

# C

## CAF

Voir le [cadre d'adoption du AWS cloud](#).

## déploiement de Canary

Diffusion lente et progressive d'une version pour les utilisateurs finaux. Lorsque vous êtes sûr, vous déployez la nouvelle version et remplacez la version actuelle dans son intégralité.

## CCoE

Voir [le Centre d'excellence du cloud](#).

## CDC

Consultez la section [Capture des données de modification](#).

## capture des données de modification (CDC)

Processus de suivi des modifications apportées à une source de données, telle qu'une table de base de données, et d'enregistrement des métadonnées relatives à ces modifications. Vous pouvez utiliser la CDC à diverses fins, telles que l'audit ou la réplication des modifications dans un système cible afin de maintenir la synchronisation.

## ingénierie du chaos

Introduire intentionnellement des défaillances ou des événements perturbateurs pour tester la résilience d'un système. Vous pouvez utiliser [AWS Fault Injection Service \(AWS FIS\)](#) pour effectuer des expériences qui stressent vos AWS charges de travail et évaluer leur réponse.

## CI/CD

Découvrez [l'intégration continue et la livraison continue](#).

## classification

Processus de catégorisation qui permet de générer des prédictions. Les modèles de ML pour les problèmes de classification prédisent une valeur discrète. Les valeurs discrètes se distinguent toujours les unes des autres. Par exemple, un modèle peut avoir besoin d'évaluer la présence ou non d'une voiture sur une image.

## Développeur citoyen

Un utilisateur professionnel qui crée des applications d'intelligence artificielle à l'aide de plateformes sans code/low code sans compétences techniques spécialisées.

## chiffrement côté client

Chiffrement des données localement, avant que la cible ne les Service AWS reçoive.

## Centre d'excellence cloud (CCoE)

Une équipe multidisciplinaire qui dirige les efforts d'adoption du cloud au sein d'une organisation, notamment en développant les bonnes pratiques en matière de cloud, en mobilisant des ressources, en établissant des délais de migration et en guidant l'organisation dans le cadre de transformations à grande échelle. Pour plus d'informations, consultez les [articles du CCoE](#) sur le blog de stratégie AWS Cloud d'entreprise.

## cloud computing

Technologie cloud généralement utilisée pour le stockage de données à distance et la gestion des appareils IoT. Le cloud computing est généralement associé à la technologie [informatique de pointe](#).

## modèle d'exploitation du cloud

Dans une organisation informatique, modèle d'exploitation utilisé pour créer, faire évoluer et optimiser un ou plusieurs environnements cloud. Pour plus d'informations, consultez la section [Création de votre modèle d'exploitation cloud](#).

## étapes d'adoption du cloud

Les quatre phases que les entreprises traversent généralement lorsqu'elles migrent vers AWS Cloud :

- **Projet** : exécution de quelques projets liés au cloud à des fins de preuve de concept et d'apprentissage
- **Base** : réaliser des investissements fondamentaux pour mettre à l'échelle l'adoption du cloud (par exemple, en créant une zone de destination, en définissant un CCoE ou en établissant un modèle opérationnel)
- **Migration** : migration d'applications individuelles
- **Re-invention** — Optimisation des produits et services et innovation dans le cloud

Ces étapes ont été définies par Stephen Orban dans le billet de blog [The Journey Toward Cloud-First & the Stages of Adoption](#) publié sur le blog AWS Cloud Enterprise Strategy. Pour plus d'informations sur leur lien avec la stratégie de AWS migration, consultez le [guide de préparation à la migration](#).

## CMDB

Consultez la base de [données de gestion des configurations](#).

## référentiel de code

Emplacement où le code source et d'autres ressources, comme la documentation, les exemples et les scripts, sont stockés et mis à jour par le biais de processus de contrôle de version. Les référentiels cloud courants incluent GitHub ou Bitbucket Cloud. Chaque version du code est appelée branche. Dans une structure de microservice, chaque référentiel est consacré à une seule fonctionnalité. Un CI/CD pipeline unique peut utiliser plusieurs référentiels.

## cache passif

Cache tampon vide, mal rempli ou contenant des données obsolètes ou non pertinentes. Cela affecte les performances, car l'instance de base de données doit lire à partir de la mémoire principale ou du disque, ce qui est plus lent que la lecture à partir du cache tampon.

## données gelées

Données rarement consultées et généralement historiques. Lorsque vous interrogez ce type de données, les requêtes lentes sont généralement acceptables. Le transfert de ces données vers des niveaux ou classes de stockage moins performants et moins coûteux peut réduire les coûts.

## vision par ordinateur (CV)

Domaine de l'[IA](#) qui utilise l'apprentissage automatique pour analyser et extraire des informations à partir de formats visuels tels que des images numériques et des vidéos. Par exemple, Amazon SageMaker AI fournit des algorithmes de traitement d'image pour les CV.

## dérive de configuration

Pour une charge de travail, une modification de configuration par rapport à l'état attendu. Cela peut entraîner une non-conformité de la charge de travail, et cela est généralement progressif et involontaire.

## base de données de gestion des configurations (CMDB)

Référentiel qui stocke et gère les informations relatives à une base de données et à son environnement informatique, y compris les composants matériels et logiciels ainsi que leurs configurations. Vous utilisez généralement les données d'une CMDB lors de la phase de découverte et d'analyse du portefeuille de la migration.

## pack de conformité

Ensemble de AWS Config règles et d'actions correctives que vous pouvez assembler pour personnaliser vos contrôles de conformité et de sécurité. Vous pouvez déployer un pack de conformité en tant qu'entité unique dans une région Compte AWS et, ou au sein d'une organisation, à l'aide d'un modèle YAML. Pour plus d'informations, consultez la section [Packs de conformité](#) dans la AWS Config documentation.

## intégration continue et livraison continue (CI/CD)

Processus d'automatisation des étapes de source, de construction, de test, de préparation et de production du processus de publication du logiciel. CI/CD est communément décrit comme un pipeline. CI/CD peut vous aider à automatiser les processus, à améliorer la productivité,

à améliorer la qualité du code et à accélérer les livraisons. Pour plus d'informations, veuillez consulter [Avantages de la livraison continue](#). CD peut également signifier déploiement continu. Pour plus d'informations, veuillez consulter [Livraison continue et déploiement continu](#).

## CV

Voir [vision par ordinateur](#).

## D

### données au repos

Données stationnaires dans votre réseau, telles que les données stockées.

### classification des données

Processus permettant d'identifier et de catégoriser les données de votre réseau en fonction de leur sévérité et de leur sensibilité. Il s'agit d'un élément essentiel de toute stratégie de gestion des risques de cybersécurité, car il vous aide à déterminer les contrôles de protection et de conservation appropriés pour les données. La classification des données est une composante du pilier de sécurité du AWS Well-Architected cadre. Pour plus d'informations, veuillez consulter [Classification des données](#).

### dérive des données

Une variation significative entre les données de production et les données utilisées pour entraîner un modèle ML, ou une modification significative des données d'entrée au fil du temps. La dérive des données peut réduire la qualité, la précision et l'équité globales des prédictions des modèles ML.

### données en transit

Données qui circulent activement sur votre réseau, par exemple entre les ressources du réseau.

### maillage de données

Un cadre architectural qui fournit une propriété des données distribuée et décentralisée avec une gestion et une gouvernance centralisées.

### minimisation des données

Le principe de collecte et de traitement des seules données strictement nécessaires. La pratique de la minimisation des données AWS Cloud peut réduire les risques liés à la confidentialité, les coûts et l'empreinte carbone de vos analyses.

## périmètre de données

Ensemble de garde-fous préventifs dans votre AWS environnement qui permettent de garantir que seules les identités fiables accèdent aux ressources fiables des réseaux attendus. Pour plus d'informations, voir [Création d'un périmètre de données sur AWS](#).

## prétraitement des données

Pour transformer les données brutes en un format facile à analyser par votre modèle de ML. Le prétraitement des données peut impliquer la suppression de certaines colonnes ou lignes et le traitement des valeurs manquantes, incohérentes ou en double.

## provenance des données

Le processus de suivi de l'origine et de l'historique des données tout au long de leur cycle de vie, par exemple la manière dont les données ont été générées, transmises et stockées.

## sujet des données

Personne dont les données sont collectées et traitées.

## entrepôt des données

Un système de gestion des données qui prend en charge les informations commerciales, telles que les analyses. Les entrepôts de données contiennent généralement de grandes quantités de données historiques et sont généralement utilisés pour les requêtes et les analyses.

## langage de définition de base de données (DDL)

Instructions ou commandes permettant de créer ou de modifier la structure des tables et des objets dans une base de données.

## langage de manipulation de base de données (DML)

Instructions ou commandes permettant de modifier (insérer, mettre à jour et supprimer) des informations dans une base de données.

## DDL

Voir [langage de définition de base](#) de données.

## ensemble profond

Sert à combiner plusieurs modèles de deep learning à des fins de prédiction. Vous pouvez utiliser des ensembles profonds pour obtenir une prévision plus précise ou pour estimer l'incertitude des prédictions.

## deep learning

Un sous-champ de ML qui utilise plusieurs couches de réseaux neuronaux artificiels pour identifier le mappage entre les données d'entrée et les variables cibles d'intérêt.

## défense en profondeur

Approche de la sécurité de l'information dans laquelle une série de mécanismes et de contrôles de sécurité sont judicieusement répartis sur l'ensemble d'un réseau informatique afin de protéger la confidentialité, l'intégrité et la disponibilité du réseau et des données qu'il contient. Lorsque vous adoptez cette stratégie AWS, vous ajoutez plusieurs contrôles à différentes couches de la AWS Organizations structure afin de sécuriser les ressources. Par exemple, une approche de défense approfondie peut combiner l'authentification multifactorielle, la segmentation du réseau et le chiffrement.

## administrateur délégué

Dans AWS Organizations, un service compatible peut enregistrer un compte AWS membre pour administrer les comptes de l'organisation et gérer les autorisations pour ce service. Ce compte est appelé administrateur délégué pour ce service. Pour plus d'informations et une liste des services compatibles, veuillez consulter la rubrique [Services qui fonctionnent avec AWS Organizations](#) dans la documentation AWS Organizations .

## déploiement

Processus de mise à disposition d'une application, de nouvelles fonctionnalités ou de corrections de code dans l'environnement cible. Le déploiement implique la mise en œuvre de modifications dans une base de code, puis la génération et l'exécution de cette base de code dans les environnements de l'application.

## environnement de développement

Voir [environnement](#).

## contrôle de détection

Contrôle de sécurité conçu pour détecter, journaliser et alerter après la survenue d'un événement. Ces contrôles constituent une deuxième ligne de défense et vous alertent en cas d'événements de sécurité qui ont contourné les contrôles préventifs en place. Pour plus d'informations, veuillez consulter la rubrique [Contrôles de détection](#) dans Implementing security controls on AWS.

## cartographie de la chaîne de valeur du développement (DVSM)

Processus utilisé pour identifier et hiérarchiser les contraintes qui nuisent à la rapidité et à la qualité du cycle de vie du développement logiciel. DVSM étend le processus de cartographie de la chaîne de valeur initialement conçu pour les pratiques de production allégée. Il met l'accent sur les étapes et les équipes nécessaires pour créer et transférer de la valeur tout au long du processus de développement logiciel.

## jumeau numérique

Représentation virtuelle d'un système réel, tel qu'un bâtiment, une usine, un équipement industriel ou une ligne de production. Les jumeaux numériques prennent en charge la maintenance prédictive, la surveillance à distance et l'optimisation de la production.

## tableau des dimensions

Dans un [schéma en étoile](#), table plus petite contenant les attributs de données relatifs aux données quantitatives d'une table de faits. Les attributs des tables de dimensions sont généralement des champs de texte ou des nombres discrets qui se comportent comme du texte. Ces attributs sont couramment utilisés pour la contrainte des requêtes, le filtrage et l'étiquetage des ensembles de résultats.

## catastrophe

Un événement qui empêche une charge de travail ou un système d'atteindre ses objectifs commerciaux sur son site de déploiement principal. Ces événements peuvent être des catastrophes naturelles, des défaillances techniques ou le résultat d'actions humaines, telles qu'une mauvaise configuration involontaire ou une attaque de logiciel malveillant.

## reprise après sinistre (DR)

La stratégie et le processus que vous utilisez pour minimiser les temps d'arrêt et les pertes de données causés par un [sinistre](#). Pour plus d'informations, consultez la section [Reprise après sinistre des charges de travail sur AWS : Restauration dans le cloud](#) dans le AWS Well-Architected Framework.

## DML

Voir [langage de manipulation de base](#) de données.

## conception axée sur le domaine

Approche visant à développer un système logiciel complexe en connectant ses composants à des domaines évolutifs, ou objectifs métier essentiels, que sert chaque composant. Ce concept

a été introduit par Eric Evans dans son livre, *Domain-Driven Design : Tackling Complexity in the Heart of Software* (Boston : Addison-Wesley Professional, 2003). Pour plus d'informations sur la manière dont vous pouvez utiliser la conception axée sur le domaine avec le modèle Strangler Fig, consultez la section [Modernisation incrémentielle des anciens services Web ASP.NET Microsoft \(ASMX\) à l'aide de conteneurs et d'Amazon API Gateway](#).

## DR

Consultez la section [Reprise après sinistre](#).

## détection de dérive

Suivi des écarts par rapport à une configuration de référence. Par exemple, vous pouvez l'utiliser AWS CloudFormation pour [détecter la dérive des ressources du système](#) ou AWS Control Tower pour [détecter les modifications de votre zone d'atterrissage](#) susceptibles d'affecter le respect des exigences de gouvernance.

## DVSM

Voir la [cartographie de la chaîne de valeur du développement](#).

## E

### EDA

Voir [analyse exploratoire des données](#).

### EDI

Voir échange [de données informatisé](#).

## informatique de périphérie

Technologie qui augmente la puissance de calcul des appareils intelligents en périphérie d'un réseau IoT. Comparé au [cloud computing, l'informatique](#) de pointe peut réduire la latence des communications et améliorer le temps de réponse.

## échange de données informatisé (EDI)

L'échange automatique de documents commerciaux entre les organisations. Pour plus d'informations, voir [Qu'est-ce que l'échange de données informatisé ?](#)

## chiffrement

Processus informatique qui transforme des données en texte clair, lisibles par l'homme, en texte chiffré.

## clé de chiffrement

Chaîne cryptographique de bits aléatoires générée par un algorithme cryptographique. La longueur des clés peut varier, et chaque clé est conçue pour être imprévisible et unique.

## endianisme

Ordre dans lequel les octets sont stockés dans la mémoire de l'ordinateur. Big-endian les systèmes stockent d'abord l'octet le plus significatif. Little-endian les systèmes stockent d'abord l'octet le moins significatif.

## point de terminaison

Voir [point de terminaison de service](#).

## service de point de terminaison

Service que vous pouvez héberger sur un cloud privé virtuel (VPC) pour le partager avec d'autres utilisateurs. Vous pouvez créer un service de point de terminaison avec AWS PrivateLink et accorder des autorisations à d'autres principaux Comptes AWS ou à Gestion des identités et des accès AWS (IAM) principaux. Ces comptes ou principaux peuvent se connecter à votre service de point de terminaison de manière privée en créant des points de terminaison d'un VPC d'interface. Pour plus d'informations, veuillez consulter [Création d'un service de point de terminaison](#) dans la documentation Amazon Virtual Private Cloud (Amazon VPC).

## planification des ressources d'entreprise (ERP)

Système qui automatise et gère les principaux processus métier (tels que la comptabilité, le [MES](#) et la gestion de projet) pour une entreprise.

## chiffrement d'enveloppe

Processus de chiffrement d'une clé de chiffrement à l'aide d'une autre clé de chiffrement. Pour plus d'informations, consultez la section [Chiffrement des enveloppes](#) dans la documentation AWS Key Management Service (AWS KMS).

## environnement

Instance d'une application en cours d'exécution. Les types d'environnement les plus courants dans le cloud computing sont les suivants :

- Environnement de développement : instance d'une application en cours d'exécution à laquelle seule l'équipe principale chargée de la maintenance de l'application peut accéder. Les environnements de développement sont utilisés pour tester les modifications avant de les promouvoir dans les environnements supérieurs. Ce type d'environnement est parfois appelé environnement de test.
- Environnements inférieurs : tous les environnements de développement d'une application, tels que ceux utilisés pour les générations et les tests initiaux.
- Environnement de production : instance d'une application en cours d'exécution à laquelle les utilisateurs finaux peuvent accéder. Dans un CI/CD pipeline, l'environnement de production est le dernier environnement de déploiement.
- Environnements supérieurs : tous les environnements accessibles aux utilisateurs autres que l'équipe de développement principale. Ils peuvent inclure un environnement de production, des environnements de préproduction et des environnements pour les tests d'acceptation par les utilisateurs.

## épopée

Dans les méthodologies agiles, catégories fonctionnelles qui aident à organiser et à prioriser votre travail. Les épopées fournissent une description détaillée des exigences et des tâches d'implémentation. Par exemple, les points forts de la AWS CAF en matière de sécurité incluent la gestion des identités et des accès, les contrôles de détection, la sécurité des infrastructures, la protection des données et la réponse aux incidents. Pour plus d'informations sur les épopées dans la stratégie de migration AWS , veuillez consulter le [guide d'implémentation du programme](#).

## ERP

Voir [Planification des ressources d'entreprise](#).

## analyse exploratoire des données (EDA)

Processus d'analyse d'un jeu de données pour comprendre ses principales caractéristiques. Vous collectez ou agrégez des données, puis vous effectuez des enquêtes initiales pour trouver des modèles, détecter des anomalies et vérifier les hypothèses. L'EDA est réalisée en calculant des statistiques récapitulatives et en créant des visualisations de données.

## F

### tableau des faits

La table centrale dans un [schéma en étoile](#). Il stocke des données quantitatives sur les opérations commerciales. Généralement, une table de faits contient deux types de colonnes : celles qui contiennent des mesures et celles qui contiennent une clé étrangère pour une table de dimensions.

### échouer rapidement

Une philosophie qui utilise des tests fréquents et progressifs pour réduire le cycle de vie du développement. C'est un élément essentiel d'une approche agile.

### limite d'isolation des défauts

Dans le AWS Cloud, une limite telle qu'une zone de disponibilité Région AWS, un plan de contrôle ou un plan de données qui limite l'effet d'une panne et contribue à améliorer la résilience des charges de travail. Pour plus d'informations, consultez la section [Limites d'isolation des AWS pannes](#).

### branche de fonctionnalités

Voir [la succursale](#).

### fonctionnalités

Les données d'entrée que vous utilisez pour faire une prédiction. Par exemple, dans un contexte de fabrication, les fonctionnalités peuvent être des images capturées périodiquement à partir de la ligne de fabrication.

### importance des fonctionnalités

Le niveau d'importance d'une fonctionnalité pour les prédictions d'un modèle. Il s'exprime généralement sous la forme d'un score numérique qui peut être calculé à l'aide de différentes techniques, telles que la méthode Shapley Additive Explanations (SHAP) et les gradients intégrés. Pour plus d'informations, voir [Interprétabilité du modèle d'apprentissage automatique avec AWS](#).

### transformation de fonctionnalité

Optimiser les données pour le processus de ML, notamment en enrichissant les données avec des sources supplémentaires, en mettant à l'échelle les valeurs ou en extrayant plusieurs ensembles d'informations à partir d'un seul champ de données. Cela permet au modèle de ML

de tirer parti des données. Par exemple, si vous décomposez la date « 2021-05-27 00:15:37 » en « 2021 », « mai », « jeudi » et « 15 », vous pouvez aider l'algorithme d'apprentissage à apprendre des modèles nuancés associés à différents composants de données.

invitation en quelques coups

Fournir à un [LLM](#) un petit nombre d'exemples illustrant la tâche et le résultat souhaité avant de lui demander d'effectuer une tâche similaire. Cette technique est une application de l'apprentissage contextuel, dans le cadre de laquelle les modèles apprennent à partir d'exemples (prises de vue) intégrés dans des instructions. Few-shot l'envoi d'instructions peut être efficace pour les tâches qui nécessitent un formatage, un raisonnement ou une connaissance du domaine spécifiques. Voir également l'[invite Zero-Shot](#).

FGAC

Découvrez le [contrôle d'accès détaillé](#).

contrôle d'accès détaillé (FGAC)

Utilisation de plusieurs conditions pour autoriser ou refuser une demande d'accès.

migration instantanée (flash-cut)

Méthode de migration de base de données qui utilise la réplication continue des données par [le biais de la capture des données de modification](#) afin de migrer les données dans les plus brefs délais, au lieu d'utiliser une approche progressive. L'objectif est de réduire au maximum les temps d'arrêt.

FM

Voir le [modèle de fondation](#).

modèle de fondation (FM)

Un vaste réseau neuronal d'apprentissage profond qui s'entraîne sur des ensembles de données massifs de données généralisées et non étiquetées. Les FM sont capables d'effectuer une grande variété de tâches générales, telles que la compréhension du langage, la génération de texte et d'images et la conversation en langage naturel. Pour plus d'informations, voir [Que sont les modèles de base ?](#)

Passerelle FM

Un intermédiaire centralisé qui contrôle et normalise l'accès aux [modèles de base](#). Également connue sous le nom de passerelle LLM.

# G

## IA générative

Sous-ensemble de modèles d'[IA](#) qui ont été entraînés sur de grandes quantités de données et qui peuvent utiliser une simple invite textuelle pour créer de nouveaux contenus et artefacts, tels que des images, des vidéos, du texte et du son. Pour plus d'informations, consultez [Qu'est-ce que l'IA générative](#).

## blocage géographique

Voir les [restrictions géographiques](#).

## restrictions géographiques (blocage géographique)

Sur Amazon CloudFront, option permettant d'empêcher les utilisateurs de certains pays d'accéder aux distributions de contenu. Vous pouvez utiliser une liste d'autorisation ou une liste de blocage pour spécifier les pays approuvés et interdits. Pour plus d'informations, consultez [la section Restreindre la distribution géographique de votre contenu](#) dans la CloudFront documentation.

## Flux de travail Gitflow

Approche dans laquelle les environnements inférieurs et supérieurs utilisent différentes branches dans un référentiel de code source. Le flux de travail Gitflow est considéré comme existant, et le [flux de travail basé sur les troncs](#) est l'approche moderne préférée.

## image dorée

Un instantané d'un système ou d'un logiciel utilisé comme modèle pour déployer de nouvelles instances de ce système ou logiciel. Par exemple, dans le secteur de la fabrication, une image dorée peut être utilisée pour fournir des logiciels sur plusieurs appareils et contribue à améliorer la vitesse, l'évolutivité et la productivité des opérations de fabrication des appareils.

## stratégie inédite

L'absence d'infrastructures existantes dans un nouvel environnement. Lorsque vous adoptez une stratégie inédite pour une architecture système, vous pouvez sélectionner toutes les nouvelles technologies sans restriction de compatibilité avec l'infrastructure existante, également appelée [brownfield](#). Si vous étendez l'infrastructure existante, vous pouvez combiner des politiques brownfield (existantes) et greenfield (inédites).

## barrière de protection

Règle de haut niveau qui permet de régir les ressources, les politiques et la conformité au sein des unités d'organisation (UO). Les barrières de protection préventives appliquent des politiques pour garantir l'alignement sur les normes de conformité. Elles sont mises en œuvre à l'aide de politiques de contrôle des services et de limites des autorisations IAM. Les barrières de protection de détection détectent les violations des politiques et les problèmes de conformité, et génèrent des alertes pour y remédier. Ils sont implémentés à l'aide d'Amazon AWS Config AWS Security Hub CSPM GuardDuty AWS Trusted Advisor, d'Amazon Inspector et de AWS Lambda contrôles personnalisés.

## rambardes (AI)

Des mécanismes de sécurité qui filtrent, valident et limitent les entrées et sorties des [agents](#) afin de garantir un comportement responsable et sûr de l'IA.

# H

## HA

Découvrez [la haute disponibilité](#).

## migration de base de données hétérogène

Migration de votre base de données source vers une base de données cible qui utilise un moteur de base de données différent (par exemple, Oracle vers Amazon Aurora). La migration hétérogène fait généralement partie d'un effort de réarchitecture, et la conversion du schéma peut s'avérer une tâche complexe. [AWS propose AWS SCT](#) qui facilite les conversions de schémas.

## haute disponibilité (HA)

Capacité d'une charge de travail à fonctionner en continu, sans intervention, en cas de difficultés ou de catastrophes. Les systèmes HA sont conçus pour basculer automatiquement, fournir constamment des performances de haute qualité et gérer différentes charges et défaillances avec un impact minimal sur les performances.

## modernisation des historiens

Approche utilisée pour moderniser et mettre à niveau les systèmes de technologie opérationnelle (OT) afin de mieux répondre aux besoins de l'industrie manufacturière. Un historien est un type

de base de données utilisé pour collecter et stocker des données provenant de diverses sources dans une usine.

#### données de rétention

Partie de données historiques étiquetées qui n'est pas divulguée dans un ensemble de données utilisé pour entraîner un modèle d'[apprentissage automatique](#). Vous pouvez utiliser les données de blocage pour évaluer les performances du modèle en comparant les prévisions du modèle aux données de blocage.

#### humain dans la boucle (HiTL)

Un modèle de flux de travail dans lequel l'exécution des [agents](#) s'arrête pour examen et approbation par l'homme aux points de décision critiques.

#### migration de base de données homogène

Migration de votre base de données source vers une base de données cible qui partage le même moteur de base de données (par exemple, Microsoft SQL Server vers Amazon RDS for SQL Server). La migration homogène s'inscrit généralement dans le cadre d'un effort de réhébergement ou de replateforme. Vous pouvez utiliser les utilitaires de base de données natifs pour migrer le schéma.

#### données chaudes

Données fréquemment consultées, telles que les données en temps réel ou les données transactionnelles récentes. Ces données nécessitent généralement un niveau ou une classe de stockage à hautes performances pour fournir des réponses rapides aux requêtes.

#### correctif

Solution d'urgence à un problème critique dans un environnement de production. En raison de son urgence, un correctif est généralement créé en dehors du flux de travail de DevOps publication habituel.

#### période de soins intensifs

Immédiatement après le basculement, période pendant laquelle une équipe de migration gère et surveille les applications migrées dans le cloud afin de résoudre les problèmes éventuels. En règle générale, cette période dure de 1 à 4 jours. À la fin de la période de soins intensifs, l'équipe de migration transfère généralement la responsabilité des applications à l'équipe des opérations cloud.

I

IaC

Considérez [l'infrastructure comme un code](#).

politique basée sur l'identité

Politique attachée à un ou plusieurs principaux IAM qui définit leurs autorisations au sein de l'AWS Cloud environnement.

application inactive

Application dont l'utilisation moyenne du processeur et de la mémoire se situe entre 5 et 20 % sur une période de 90 jours. Dans un projet de migration, il est courant de retirer ces applications ou de les retenir sur site.

IIoT

Voir [Internet industriel des objets](#).

infrastructure immuable

Modèle qui déploie une nouvelle infrastructure pour les charges de travail de production au lieu de mettre à jour, d'appliquer des correctifs ou de modifier l'infrastructure existante. Les infrastructures immuables sont intrinsèquement plus cohérentes, fiables et prévisibles que les infrastructures [mutables](#). Pour plus d'informations, consultez les meilleures pratiques de [déploiement à l'aide d'une infrastructure immuable](#) dans le AWS Well-Architected Framework.

VPC entrant (d'entrée)

Dans une architecture AWS multi-comptes, un VPC qui accepte, inspecte et achemine les connexions réseau depuis l'extérieur d'une application. L'[architecture de référence de sécuritéAWS](#) recommande de configurer votre compte réseau avec des VPC entrants, sortants et d'inspection afin de protéger l'interface bidirectionnelle entre votre application et Internet en général.

migration incrémentielle

Stratégie de basculement dans le cadre de laquelle vous migrez votre application par petites parties au lieu d'effectuer un basculement complet unique. Par exemple, il se peut que vous ne transfériez que quelques microservices ou utilisateurs vers le nouveau système dans un

I

premier temps. Après avoir vérifié que tout fonctionne correctement, vous pouvez transférer progressivement des microservices ou des utilisateurs supplémentaires jusqu'à ce que vous puissiez mettre hors service votre système hérité. Cette stratégie réduit les risques associés aux migrations de grande ampleur.

## Industry 4.0

Terme introduit par [Klaus Schwab](#) en 2016 pour désigner la modernisation des processus de fabrication grâce aux avancées en matière de connectivité, de données en temps réel, d'automatisation, d'analyse et. AI/ML

## infrastructure

Ensemble des ressources et des actifs contenus dans l'environnement d'une application.

## infrastructure en tant que code (IaC)

Processus de mise en service et de gestion de l'infrastructure d'une application via un ensemble de fichiers de configuration. IaC est conçue pour vous aider à centraliser la gestion de l'infrastructure, à normaliser les ressources et à mettre à l'échelle rapidement afin que les nouveaux environnements soient reproductibles, fiables et cohérents.

## internet industriel des objets (IIoT)

L'utilisation de capteurs et d'appareils connectés à Internet dans les secteurs industriels tels que la fabrication, l'énergie, l'automobile, les soins de santé, les sciences de la vie et l'agriculture. Pour plus d'informations, veuillez consulter [Building an industrial Internet of Things \(IIoT\) digital transformation strategy](#).

## VPC d'inspection

Dans une architecture AWS multi-comptes, un VPC centralisé qui gère les inspections du trafic réseau entre les VPC (identiques ou Régions AWS différents), Internet et les réseaux sur site. L'[architecture de référence de sécuritéAWS](#) recommande de configurer votre compte réseau avec des VPC entrants, sortants et d'inspection afin de protéger l'interface bidirectionnelle entre votre application et Internet en général.

## Internet des objets (IoT)

Réseau d'objets physiques connectés dotés de capteurs ou de processeurs intégrés qui communiquent avec d'autres appareils et systèmes via Internet ou via un réseau de communication local. Pour plus d'informations, veuillez consulter la section [Qu'est-ce que l'IoT ?](#).

## interprétabilité

Caractéristique d'un modèle de machine learning qui décrit dans quelle mesure un être humain peut comprendre comment les prédictions du modèle dépendent de ses entrées. Pour plus d'informations, voir [Interprétabilité du modèle d'apprentissage automatique avec AWS](#).

## IoT

Voir [Internet des objets](#).

## Bibliothèque d'informations informatiques (ITIL)

Ensemble de bonnes pratiques pour proposer des services informatiques et les aligner sur les exigences métier. L'ITIL constitue la base de l'ITSM.

## gestion des services informatiques (ITSM)

Activités associées à la conception, à la mise en œuvre, à la gestion et à la prise en charge de services informatiques d'une organisation. Pour plus d'informations sur l'intégration des opérations cloud aux outils ITSM, veuillez consulter le [guide d'intégration des opérations](#).

## ITIL

Consultez la [bibliothèque d'informations informatiques](#).

## ITSM

Voir [Gestion des services informatiques](#).

## L

### contrôle d'accès basé sur des étiquettes (LBAC)

Une implémentation du contrôle d'accès obligatoire (MAC) dans laquelle une valeur d'étiquette de sécurité est explicitement attribuée aux utilisateurs et aux données elles-mêmes. L'intersection entre l'étiquette de sécurité utilisateur et l'étiquette de sécurité des données détermine les lignes et les colonnes visibles par l'utilisateur.

### zone de destination

Une zone d'atterrissage est un AWS environnement multi-comptes bien conçu, évolutif et sécurisé. Il s'agit d'un point de départ à partir duquel vos entreprises peuvent rapidement lancer et déployer des charges de travail et des applications en toute confiance dans leur environnement

de sécurité et d'infrastructure. Pour plus d'informations sur les zones de destination, veuillez consulter [Setting up a secure and scalable multi-account AWS environment](#).

grand modèle de langage (LLM)

Un modèle d'[intelligence artificielle basé](#) sur le deep learning qui est préentraîné sur une grande quantité de données. Un LLM peut effectuer plusieurs tâches, telles que répondre à des questions, résumer des documents, traduire du texte dans d'autres langues et compléter des phrases. Pour plus d'informations, voir [Que sont les LLM](#).

migration de grande envergure

Migration de 300 serveurs ou plus.

LBAC

Voir contrôle d'[accès basé sur des étiquettes](#).

principe de moindre privilège

Bonne pratique de sécurité qui consiste à accorder les autorisations minimales nécessaires à l'exécution d'une tâche. Pour plus d'informations, veuillez consulter la rubrique [Accorder les autorisations de moindre privilège](#) dans la documentation IAM.

lift and shift

Voir [7 Rs](#).

système de poids faible

Système qui stocke d'abord l'octet le moins significatif. Voir aussi [endianité](#).

LLM

Voir le [grand modèle de langage](#).

environnements inférieurs

Voir [environnement](#).

## M

machine learning (ML)

Type d'intelligence artificielle qui utilise des algorithmes et des techniques pour la reconnaissance et l'apprentissage de modèles. Le ML analyse et apprend à partir de données enregistrées, telles

que les données de l'Internet des objets (IoT), pour générer un modèle statistique basé sur des modèles. Pour plus d'informations, veuillez consulter [Machine Learning](#).

branche principale

Voir [la succursale](#).

malware

Logiciel conçu pour compromettre la sécurité ou la confidentialité de l'ordinateur. Les logiciels malveillants peuvent perturber les systèmes informatiques, divulguer des informations sensibles ou obtenir un accès non autorisé. Parmi les malwares, on peut citer les virus, les vers, les rançongiciels, les chevaux de Troie, les logiciels espions et les enregistreurs de frappe.

services gérés

Services AWS pour lequel AWS fonctionnent la couche d'infrastructure, le système d'exploitation et les plateformes, et vous accédez aux points de terminaison pour stocker et récupérer des données. Amazon Simple Storage Service (Amazon S3) et Amazon DynamoDB sont des exemples de services gérés. Ils sont également connus sous le nom de services abstraits.

système d'exécution de la fabrication (MES)

Un système logiciel pour le suivi, la surveillance, la documentation et le contrôle des processus de production qui convertissent les matières premières en produits finis dans l'atelier.

MAP

Voir [Migration Acceleration Program](#).

MCP

Voir [Model Context Protocol](#).

Protocole de contexte du modèle (MCP)

Protocole sans état pour la communication entre [un agent](#) et un [outil](#).

serveur MCP

Service qui expose un ou plusieurs [outils](#) via le [protocole Model Context](#).

mécanisme

Processus complet au cours duquel vous créez un outil, favorisez son adoption, puis inspectez les résultats afin de procéder aux ajustements nécessaires. Un mécanisme est un cycle qui se

renforce et s'améliore au fur et à mesure de son fonctionnement. Pour plus d'informations, voir [Création de mécanismes](#) dans le AWS Well-Architected cadre.

## compte membre

Tous, à l'exception des Comptes AWS exception du compte de gestion, qui font partie d'une organisation dans AWS Organizations. Un compte ne peut être membre que d'une seule organisation à la fois.

## MAILLES

Voir le [système d'exécution de la fabrication](#).

## Transport télémétrique en file d'attente de messages (MQTT)

[Un protocole de communication léger de machine à machine \(M2M\), basé sur le publish/subscribe modèle, pour les appareils IoT aux ressources limitées.](#)

## microservice

Petit service indépendant qui communique via des API bien définies et qui est généralement détenu par de petites équipes autonomes. Par exemple, un système d'assurance peut inclure des microservices qui mappent à des capacités métier, telles que les ventes ou le marketing, ou à des sous-domaines, tels que les achats, les réclamations ou l'analytique. Les avantages des microservices incluent l'agilité, la flexibilité de la mise à l'échelle, la facilité de déploiement, la réutilisation du code et la résilience. Pour plus d'informations, consultez la section [Intégration de microservices à l'aide de services AWS sans serveur](#).

## architecture de microservices

Approche de création d'une application avec des composants indépendants qui exécutent chaque processus d'application en tant que microservice. Ces microservices communiquent via une interface bien définie à l'aide d'API légères. Chaque microservice de cette architecture peut être mis à jour, déployé et mis à l'échelle pour répondre à la demande de fonctions spécifiques d'une application. Pour plus d'informations, consultez la section [Implémentation de microservices sur AWS](#).

## Programme d'accélération des migrations (MAP)

Un AWS programme qui fournit un support de conseil, des formations et des services pour aider les entreprises à établir une base opérationnelle solide pour passer au cloud, et pour aider à compenser le coût initial des migrations. MAP inclut une méthodologie de migration pour exécuter les migrations héritées de manière méthodique, ainsi qu'un ensemble d'outils pour automatiser et accélérer les scénarios de migration courants.

## migration à grande échelle

Processus consistant à transférer la majeure partie du portefeuille d'applications vers le cloud par vagues, un plus grand nombre d'applications étant déplacées plus rapidement à chaque vague. Cette phase utilise les bonnes pratiques et les enseignements tirés des phases précédentes pour implémenter une usine de migration d'équipes, d'outils et de processus en vue de rationaliser la migration des charges de travail grâce à l'automatisation et à la livraison agile. Il s'agit de la troisième phase de la [stratégie de migration AWS](#).

## usine de migration

Cross-functional des équipes qui rationalisent la migration des charges de travail grâce à des approches automatisées et agiles. Les équipes de Migration Factory comprennent généralement des responsables des opérations, des analystes commerciaux et des propriétaires, des ingénieurs de migration, des développeurs et DevOps des professionnels travaillant dans le cadre de sprints. Entre 20 et 50 % du portefeuille d'applications d'entreprise est constitué de modèles répétés qui peuvent être optimisés par une approche d'usine. Pour plus d'informations, veuillez consulter la rubrique [discussion of migration factories](#) et le [guide Cloud Migration Factory](#) dans cet ensemble de contenus.

## métadonnées de migration

Informations relatives à l'application et au serveur nécessaires pour finaliser la migration. Chaque modèle de migration nécessite un ensemble de métadonnées de migration différent. Les exemples de métadonnées de migration incluent le sous-réseau cible, le groupe de sécurité et le AWS compte.

## modèle de migration

Tâche de migration reproductible qui détaille la stratégie de migration, la destination de la migration et l'application ou le service de migration utilisé. Exemple : réorganisez la migration vers Amazon EC2 AWS avec le service de migration d'applications.

## Évaluation du portefeuille de migration (MPA)

Outil en ligne qui fournit des informations pour valider l'analyse de rentabilisation en faveur de la migration vers le. AWS Cloud La MPA propose une évaluation détaillée du portefeuille (dimensionnement approprié des serveurs, tarification, comparaison du coût total de possession, analyse des coûts de migration), ainsi que la planification de la migration (analyse et collecte des données d'applications, regroupement des applications, priorisation des migrations et planification des vagues). L'[outil MPA](#) (connexion requise) est disponible gratuitement pour tous les AWS consultants et consultants APN Partner.

## Évaluation de la préparation à la migration (MRA)

Processus qui consiste à obtenir des informations sur l'état de préparation d'une organisation au cloud, à identifier les forces et les faiblesses et à élaborer un plan d'action pour combler les lacunes identifiées, à l'aide du AWS CAF. Pour plus d'informations, veuillez consulter le [guide de préparation à la migration](#). La MRA est la première phase de la [stratégie de migration AWS](#).

### stratégie de migration

L'approche utilisée pour migrer une charge de travail vers le AWS Cloud. Pour plus d'informations, reportez-vous aux [7 R](#) de ce glossaire et à [Mobiliser votre organisation pour accélérer les migrations à grande échelle](#).

### ML

Voir [apprentissage automatique](#).

### modernisation

Transformation d'une application obsolète (héritée ou monolithique) et de son infrastructure en un système agile, élastique et hautement disponible dans le cloud afin de réduire les coûts, de gagner en efficacité et de tirer parti des innovations. Pour plus d'informations, consultez [la section Stratégie de modernisation des applications dans le AWS Cloud](#).

### évaluation de la préparation à la modernisation

Évaluation qui permet de déterminer si les applications d'une organisation sont prêtes à être modernisées, d'identifier les avantages, les risques et les dépendances, et qui détermine dans quelle mesure l'organisation peut prendre en charge l'état futur de ces applications. Le résultat de l'évaluation est un plan de l'architecture cible, une feuille de route détaillant les phases de développement et les étapes du processus de modernisation, ainsi qu'un plan d'action pour combler les lacunes identifiées. Pour plus d'informations, consultez la section [Évaluation de l'état de préparation à la modernisation des applications dans le AWS Cloud](#).

### applications monolithiques (monolithes)

Applications qui s'exécutent en tant que service unique avec des processus étroitement couplés. Les applications monolithiques ont plusieurs inconvénients. Si une fonctionnalité de l'application connaît un pic de demande, l'architecture entière doit être mise à l'échelle. L'ajout ou l'amélioration des fonctionnalités d'une application monolithique devient également plus complexe lorsque la base de code s'élargit. Pour résoudre ces problèmes, vous pouvez utiliser une architecture de microservices. Pour plus d'informations, veuillez consulter [Decomposing monoliths into microservices](#).

## MPA

Voir [Évaluation du portefeuille de migration](#).

## MQTT

Voir [Message Queuing Telemetry Transport](#).

## classification multi-classes

Processus qui permet de générer des prédictions pour plusieurs classes (prédiction d'un résultat parmi plus de deux). Par exemple, un modèle de ML peut demander « Ce produit est-il un livre, une voiture ou un téléphone ? » ou « Quelle catégorie de produits intéresse le plus ce client ? ».

## infrastructure mutable

Modèle qui met à jour et modifie l'infrastructure existante pour les charges de travail de production. Pour améliorer la cohérence, la fiabilité et la prévisibilité, le AWS Well-Architected Framework recommande l'utilisation d'une [infrastructure immuable](#) comme meilleure pratique.

## O

### OAC

Voir [Contrôle d'accès à l'origine](#).

### OAI

Voir [l'identité d'accès à l'origine](#).

### OCM

Voir [gestion du changement organisationnel](#).

## migration hors ligne

Méthode de migration dans laquelle la charge de travail source est supprimée au cours du processus de migration. Cette méthode implique un temps d'arrêt prolongé et est généralement utilisée pour de petites charges de travail non critiques.

## OI

Consultez la section [Intégration des opérations](#).

## OLA

Voir l'accord [au niveau opérationnel](#).

## migration en ligne

Méthode de migration dans laquelle la charge de travail source est copiée sur le système cible sans être mise hors ligne. Les applications connectées à la charge de travail peuvent continuer à fonctionner pendant la migration. Cette méthode implique un temps d'arrêt nul ou minimal et est généralement utilisée pour les charges de travail de production critiques.

## OPC-UA

Voir [Open Process Communications - Architecture unifiée](#).

## Communications par processus ouvert - Architecture unifiée (OPC-UA)

Protocole de communication machine à machine (M2M) pour l'automatisation industrielle. OPC-UA fournit une norme d'interopérabilité avec des schémas de chiffrement, d'authentification et d'autorisation des données.

## accord au niveau opérationnel (OLA)

Accord qui précise ce que les groupes informatiques fonctionnels s'engagent à fournir les uns aux autres, afin de prendre en charge un contrat de niveau de service (SLA).

## examen de l'état de préparation opérationnelle (ORR)

Une liste de questions et de bonnes pratiques associées qui vous aident à comprendre, à évaluer, à prévenir ou à réduire l'ampleur des incidents et des défaillances possibles. Pour plus d'informations, voir [Examens de l'état de préparation opérationnelle \(ORR\)](#) dans le AWS Well-Architected cadre.

## technologie opérationnelle (OT)

Systèmes matériels et logiciels qui fonctionnent avec l'environnement physique pour contrôler les opérations, les équipements et les infrastructures industriels. Dans le secteur manufacturier, l'intégration des systèmes OT et des technologies de l'information (IT) est au cœur des transformations de [l'industrie 4.0](#).

## intégration des opérations (OI)

Processus de modernisation des opérations dans le cloud, qui implique la planification de la préparation, l'automatisation et l'intégration. Pour en savoir plus, veuillez consulter le [guide d'intégration des opérations](#).

## journal de suivi d'organisation

Un parcours créé par AWS CloudTrail qui enregistre tous les événements pour tous les membres Comptes AWS d'une organisation dans AWS Organizations. Ce journal de suivi est créé dans chaque Compte AWS qui fait partie de l'organisation et suit l'activité de chaque compte. Pour plus d'informations, consultez [la section Création d'un suivi pour une organisation](#) dans la CloudTrail documentation.

## gestion du changement organisationnel (OCM)

Cadre pour gérer les transformations métier majeures et perturbatrices du point de vue des personnes, de la culture et du leadership. L'OCM aide les organisations à se préparer et à effectuer la transition vers de nouveaux systèmes et de nouvelles politiques en accélérant l'adoption des changements, en abordant les problèmes de transition et en favorisant des changements culturels et organisationnels. Dans la stratégie de AWS migration, ce cadre est appelé accélération du personnel, en raison de la rapidité du changement requise dans les projets d'adoption du cloud. Pour plus d'informations, veuillez consulter le [guide OCM](#).

## contrôle d'accès d'origine (OAC)

Dans CloudFront, une option améliorée pour restreindre l'accès afin de sécuriser votre contenu Amazon Simple Storage Service (Amazon S3). OAC prend en charge tous les compartiments S3 dans leur ensemble Régions AWS, le chiffrement côté serveur avec AWS KMS (SSE-KMS) et les DELETE requêtes dynamiques PUT adressées au compartiment S3.

## identité d'accès d'origine (OAI)

Dans CloudFront, une option permettant de restreindre l'accès afin de sécuriser votre contenu Amazon S3. Lorsque vous utilisez OAI, il CloudFront crée un principal auprès duquel Amazon S3 peut s'authentifier. Les principaux authentifiés ne peuvent accéder au contenu d'un compartiment S3 que par le biais d'une distribution spécifique CloudFront . Voir également [OAC](#), qui fournit un contrôle d'accès plus précis et amélioré.

## ORR

Voir l'[examen de l'état de préparation opérationnelle](#).

## DE

Voir [technologie opérationnelle](#).

## VPC sortant (de sortie)

Dans une architecture AWS multi-comptes, un VPC qui gère les connexions réseau initiées depuis une application. L'[architecture de référence de sécurité AWS](#) recommande de configurer votre compte réseau avec des VPC entrants, sortants et d'inspection afin de protéger l'interface bidirectionnelle entre votre application et Internet en général.

## P

### limite des autorisations

Politique de gestion IAM attachée aux principaux IAM pour définir les autorisations maximales que peut avoir l'utilisateur ou le rôle. Pour plus d'informations, veuillez consulter la rubrique [Limites des autorisations](#) dans la documentation IAM.

### informations personnelles identifiables (PII)

Informations qui, lorsqu'elles sont consultées directement ou associées à d'autres données connexes, peuvent être utilisées pour déduire raisonnablement l'identité d'une personne. Les exemples d'informations personnelles incluent les noms, les adresses et les informations de contact.

### PII

Voir les [informations personnelles identifiables](#).

### manuel stratégique

Ensemble d'étapes prédéfinies qui capturent le travail associé aux migrations, comme la fourniture de fonctions d'opérations de base dans le cloud. Un manuel stratégique peut revêtir la forme de scripts, de runbooks automatisés ou d'un résumé des processus ou des étapes nécessaires au fonctionnement de votre environnement modernisé.

### PLC

Voir [contrôleur logique programmable](#).

### PLM

Consultez la section [Gestion du cycle de vie des produits](#).

## policy

Objet capable de définir les autorisations (voir la [politique basée sur l'identité](#)), de spécifier les conditions d'accès (voir la [politique basée sur les ressources](#)) ou de définir les autorisations maximales pour tous les comptes d'une organisation dans AWS Organizations (voir la politique de contrôle des [services](#)).

## persistance polyglotte

Choix indépendant de la technologie de stockage de données d'un microservice en fonction des modèles d'accès aux données et d'autres exigences. Si vos microservices utilisent la même technologie de stockage de données, ils peuvent rencontrer des difficultés d'implémentation ou présenter des performances médiocres. Les microservices sont plus faciles à mettre en œuvre, atteignent de meilleures performances, ainsi qu'une meilleure capacité de mise à l'échelle s'ils utilisent l'entrepôt de données le mieux adapté à leurs besoins.

## évaluation du portefeuille

Processus de découverte, d'analyse et de priorisation du portefeuille d'applications afin de planifier la migration. Pour plus d'informations, veuillez consulter [Evaluating migration readiness](#).

## predicate

Une condition de requête qui renvoie `true` ou `false`, généralement située dans une `WHERE` clause.

## prédicat pushdown

Technique d'optimisation des requêtes de base de données qui filtre les données de la requête avant le transfert. Cela réduit la quantité de données qui doivent être extraites et traitées à partir de la base de données relationnelle et améliore les performances des requêtes.

## contrôle préventif

Contrôle de sécurité conçu pour empêcher qu'un événement ne se produise. Ces contrôles constituent une première ligne de défense pour empêcher tout accès non autorisé ou toute modification indésirable de votre réseau. Pour plus d'informations, veuillez consulter [Preventative controls](#) dans *Implementing security controls on AWS*.

## principal

Entité capable d'effectuer AWS des actions et d'accéder à des ressources. Cette entité est généralement un utilisateur root pour un Compte AWS rôle IAM ou un utilisateur. Pour plus

d'informations, veuillez consulter la rubrique Principal dans [Termes et concepts relatifs aux rôles](#), dans la documentation IAM.

#### confidentialité dès la conception

Une approche d'ingénierie système qui prend en compte la confidentialité tout au long du processus de développement.

#### zones hébergées privées

Conteneur qui contient des informations concernant la façon dont vous souhaitez qu'Amazon Route 53 réponde aux requêtes DNS pour un domaine et ses sous-domaines dans un ou plusieurs VPC. Pour plus d'informations, veuillez consulter [Working with private hosted zones](#) dans la documentation Route 53.

#### contrôle proactif

[Contrôle de sécurité](#) conçu pour empêcher le déploiement de ressources non conformes. Ces contrôles analysent les ressources avant qu'elles ne soient provisionnées. Si la ressource n'est pas conforme au contrôle, elle n'est pas provisionnée. Pour plus d'informations, consultez le [guide de référence sur les contrôles](#) dans la AWS Control Tower documentation et consultez la section [Contrôles proactifs dans Implémentation](#) des contrôles de sécurité sur AWS.

#### gestion du cycle de vie des produits (PLM)

Gestion des données et des processus d'un produit tout au long de son cycle de vie, depuis la conception, le développement et le lancement, en passant par la croissance et la maturité, jusqu'au déclin et au retrait.

#### environnement de production

Voir [environnement](#).

#### contrôleur logique programmable (PLC)

Dans le secteur manufacturier, un ordinateur hautement fiable et adaptable qui surveille les machines et automatise les processus de fabrication.

#### chaînage rapide

Utiliser le résultat d'une invite [LLM](#) comme entrée pour l'invite suivante afin de générer de meilleures réponses. Cette technique est utilisée pour décomposer une tâche complexe en sous-tâches ou pour affiner ou développer de manière itérative une réponse préliminaire. Cela permet d'améliorer la précision et la pertinence des réponses d'un modèle et permet d'obtenir des résultats plus précis et personnalisés.

## pseudonymisation

Processus de remplacement des identifiants personnels dans un ensemble de données par des valeurs fictives. La pseudonymisation peut contribuer à protéger la vie privée. Les données pseudonymisées sont toujours considérées comme des données personnelles.

## publish/subscribe (pub/sub)

Modèle qui permet des communications asynchrones entre les microservices afin d'améliorer l'évolutivité et la réactivité. Par exemple, dans un [MES](#) basé sur des microservices, un microservice peut publier des messages d'événements sur un canal auquel d'autres microservices peuvent s'abonner. Le système peut ajouter de nouveaux microservices sans modifier le service de publication.

## Q

### plan de requête

Série d'étapes, telles que des instructions, utilisées pour accéder aux données d'un système de base de données relationnelle SQL.

### régression du plan de requêtes

Le cas où un optimiseur de service de base de données choisit un plan moins optimal qu'avant une modification donnée de l'environnement de base de données. Cela peut être dû à des changements en termes de statistiques, de contraintes, de paramètres d'environnement, de liaisons de paramètres de requêtes et de mises à jour du moteur de base de données.

## R

### Matrice RACI

Voir [responsable, responsable, consulté, informé \(RACI\)](#).

### RAG

Voir [Retrieval Augmented Generation](#).

### rançongiciel

Logiciel malveillant conçu pour bloquer l'accès à un système informatique ou à des données jusqu'à ce qu'un paiement soit effectué.

## Matrice RASCI

Voir [responsable, responsable, consulté, informé \(RACI\)](#).

## RCAC

Voir [contrôle d'accès aux lignes et aux colonnes](#).

## réplica en lecture

Copie d'une base de données utilisée en lecture seule. Vous pouvez acheminer les requêtes vers le réplica de lecture pour réduire la charge sur votre base de données principale.

## réarchitecte

Voir [7 Rs](#).

## objectif de point de récupération (RPO)

Durée maximale acceptable depuis le dernier point de récupération des données. Il détermine ce qui est considéré comme étant une perte de données acceptable entre le dernier point de reprise et l'interruption du service.

## objectif de temps de récupération (RTO)

Le délai maximum acceptable entre l'interruption du service et le rétablissement du service.

## refactoriser

Voir [7 Rs](#).

## Région

Un ensemble de AWS ressources dans une zone géographique. Chacune Région AWS est isolée et indépendante des autres pour garantir tolérance aux pannes, stabilité et résilience. Pour plus d'informations, voir [Spécifier ce que Régions AWS votre compte peut utiliser](#).

## régression

Technique de ML qui prédit une valeur numérique. Par exemple, pour résoudre le problème « Quel sera le prix de vente de cette maison ? », un modèle de ML pourrait utiliser un modèle de régression linéaire pour prédire le prix de vente d'une maison sur la base de faits connus à son sujet (par exemple, la superficie en mètres carrés).

## réhéberger

Voir [7 Rs](#).

## version

Dans un processus de déploiement, action visant à promouvoir les modifications apportées à un environnement de production.

## déplacer

Voir [7 Rs](#).

## replateforme

Voir [7 Rs](#).

## rachat

Voir [7 Rs](#).

## résilience

La capacité d'une application à résister aux perturbations ou à s'en remettre. [La haute disponibilité et la reprise après sinistre](#) sont des considérations courantes lors de la planification de la résilience dans le AWS Cloud. Pour plus d'informations, consultez la section [AWS Cloud Résilience](#).

## politique basée sur les ressources

Politique attachée à une ressource, comme un compartiment Amazon S3, un point de terminaison ou une clé de chiffrement. Ce type de politique précise les principaux auxquels l'accès est autorisé, les actions prises en charge et toutes les autres conditions qui doivent être remplies.

## matrice responsable, redevable, consulté et informé (RACI)

Une matrice qui définit les rôles et les responsabilités de toutes les parties impliquées dans les activités de migration et les opérations cloud. Le nom de la matrice est dérivé des types de responsabilité définis dans la matrice : responsable (R), responsable (A), consulté (C) et informé (I). Le type de support (S) est facultatif. Si vous incluez le support, la matrice est appelée matrice RASCI, et si vous l'excluez, elle est appelée matrice RACI.

## contrôle réactif

Contrôle de sécurité conçu pour permettre de remédier aux événements indésirables ou aux écarts par rapport à votre référence de sécurité. Pour plus d'informations, veuillez consulter la rubrique [Responsive controls](#) dans Implementing security controls on AWS.

## retain

Voir [7 Rs](#).

se retirer

Voir [7 Rs](#).

Génération augmentée de récupération (RAG)

Technologie d'[IA générative](#) dans laquelle un [LLM](#) fait référence à une source de données faisant autorité qui se trouve en dehors de ses sources de données de formation avant de générer une réponse. Par exemple, un modèle RAG peut effectuer une recherche sémantique dans la base de connaissances ou dans les données personnalisées d'une organisation. Pour plus d'informations, voir [Qu'est-ce que RAG ?](#)

rotation

Processus de mise à jour périodique d'un [secret](#) pour empêcher un attaquant d'accéder aux informations d'identification.

contrôle d'accès aux lignes et aux colonnes (RCAC)

Utilisation d'expressions SQL simples et flexibles dotées de règles d'accès définies. Le RCAC comprend des autorisations de ligne et des masques de colonnes.

RPO

Voir l'[objectif du point de récupération](#).

RTO

Voir l'[objectif en matière de temps de rétablissement](#).

runbook

Ensemble de procédures manuelles ou automatisées nécessaires à l'exécution d'une tâche spécifique. Elles visent généralement à rationaliser les opérations ou les procédures répétitives présentant des taux d'erreur élevés.

## S

SAML 2.0

Un standard ouvert utilisé par de nombreux fournisseurs d'identité (IdPs). Cette fonctionnalité permet l'authentification unique fédérée (SSO), afin que les utilisateurs puissent se connecter

AWS Management Console ou appeler les opérations de l' AWS API sans que vous ayez à créer un utilisateur dans IAM pour tous les membres de votre organisation. Pour plus d'informations sur la fédération SAML 2.0, veuillez consulter [À propos de la fédération SAML 2.0](#) dans la documentation IAM.

## SCADA

Voir [Contrôle de supervision et acquisition de données](#).

## SCP

Voir la [politique de contrôle des services](#).

## secret

Dans AWS Secrets Manager des informations confidentielles ou restreintes, telles qu'un mot de passe ou des informations d'identification utilisateur, que vous stockez sous forme cryptée. Il comprend la valeur secrète et ses métadonnées. La valeur secrète peut être binaire, une chaîne unique ou plusieurs chaînes. Pour plus d'informations, voir [Que contient le secret d'un Secrets Manager ?](#) dans la documentation de Secrets Manager.

## sécurité dès la conception

Une approche d'ingénierie système qui prend en compte la sécurité tout au long du processus de développement.

## contrôle de sécurité

Barrière de protection technique ou administrative qui empêche, détecte ou réduit la capacité d'un assaillant d'exploiter une vulnérabilité de sécurité. Il existe quatre principaux types de contrôles de sécurité : [préventifs](#), [détectifs](#), [réactifs](#) et [proactifs](#).

## renforcement de la sécurité

Processus qui consiste à réduire la surface d'attaque pour la rendre plus résistante aux attaques. Cela peut inclure des actions telles que la suppression de ressources qui ne sont plus requises, la mise en œuvre des bonnes pratiques de sécurité consistant à accorder le moindre privilège ou la désactivation de fonctionnalités inutiles dans les fichiers de configuration.

## système de gestion des informations et des événements de sécurité (SIEM)

Outils et services qui associent les systèmes de gestion des informations de sécurité (SIM) et de gestion des événements de sécurité (SEM). Un système SIEM collecte, surveille et analyse les

données provenant de serveurs, de réseaux, d'appareils et d'autres sources afin de détecter les menaces et les failles de sécurité, mais aussi de générer des alertes.

#### automatisation des réponses de sécurité

Action prédéfinie et programmée conçue pour répondre automatiquement à un événement de sécurité ou y remédier. Ces automatisations servent de contrôles de sécurité [détectifs ou réactifs](#) qui vous aident à mettre en œuvre les meilleures pratiques en matière AWS de sécurité. Parmi les actions de réponse automatique, citons la modification d'un groupe de sécurité VPC, l'application de correctifs à une instance Amazon EC2 ou la rotation des informations d'identification.

#### chiffrement côté serveur

Chiffrement des données à destination, par celui Service AWS qui les reçoit.

#### Politique de contrôle des services (SCP)

Politique qui propose un contrôle centralisé des autorisations pour tous les comptes d'une organisation dans AWS Organizations. Les SCP définissent des barrières de protection ou des limites aux actions qu'un administrateur peut déléguer à des utilisateurs ou à des rôles. Vous pouvez utiliser les SCP comme listes d'autorisation ou de refus, pour indiquer les services ou les actions autorisés ou interdits. Pour plus d'informations, consultez la section [Politiques de contrôle des services](#) dans la AWS Organizations documentation.

#### point de terminaison du service

URL du point d'entrée pour un Service AWS. Pour vous connecter par programmation au service cible, vous pouvez utiliser un point de terminaison. Pour plus d'informations, veuillez consulter la rubrique [Service AWS endpoints](#) dans Références générales AWS.

#### contrat de niveau de service (SLA)

Accord qui précise ce qu'une équipe informatique promet de fournir à ses clients, comme le temps de disponibilité et les performances des services.

#### indicateur de niveau de service (SLI)

Mesure d'un aspect des performances d'un service, tel que son taux d'erreur, sa disponibilité ou son débit.

#### objectif de niveau de service (SLO)

Mesure cible qui représente l'état d'un service, tel que mesuré par un indicateur de [niveau de service](#).

## modèle de responsabilité partagée

Un modèle décrivant la responsabilité que vous partagez en matière AWS de sécurité et de conformité dans le cloud. AWS est responsable de la sécurité du cloud, alors que vous êtes responsable de la sécurité dans le cloud. Pour de plus amples informations, veuillez consulter [Modèle de responsabilité partagée](#).

## IA de l'ombre

Applications d'[IA](#) non autorisées créées ou utilisées en dehors des canaux régis au sein d'une organisation.

## SIEM

Consultez les [informations de sécurité et le système de gestion des événements](#).

## point de défaillance unique (SPOF)

Défaillance d'un seul composant critique d'une application susceptible de perturber le système.

## SLA

Voir le contrat [de niveau de service](#).

## SLI

Voir l'indicateur de [niveau de service](#).

## SLO

Voir l'objectif de [niveau de service](#).

## modèle split-and-seed

Modèle permettant de mettre à l'échelle et d'accélérer les projets de modernisation. Au fur et à mesure que les nouvelles fonctionnalités et les nouvelles versions de produits sont définies, l'équipe principale se divise pour créer des équipes de produit. Cela permet de mettre à l'échelle les capacités et les services de votre organisation, d'améliorer la productivité des développeurs et de favoriser une innovation rapide. Pour plus d'informations, consultez la section [Approche progressive de la modernisation des applications dans le AWS Cloud](#)

## SPOF

Voir [point de défaillance unique](#).

## schéma en étoile

Structure organisationnelle de base de données qui utilise une grande table de faits pour stocker les données transactionnelles ou mesurées et utilise une ou plusieurs tables dimensionnelles plus petites pour stocker les attributs des données. Cette structure est conçue pour être utilisée dans un [entrepôt de données](#) ou à des fins de business intelligence.

## modèle de figuier étrangleur

Approche de modernisation des systèmes monolithiques en réécrivant et en remplaçant progressivement les fonctionnalités du système jusqu'à ce que le système hérité puisse être mis hors service. Ce modèle utilise l'analogie d'un figuier de vigne qui se développe dans un arbre existant et qui finit par supplanter son hôte. Le schéma a été [présenté par Martin Fowler](#) comme un moyen de gérer les risques lors de la réécriture de systèmes monolithiques. Pour un exemple d'application de ce modèle, consultez la section [Modernisation progressive des anciens services Web Microsoft ASP.NET \(ASMX\) à l'aide de conteneurs et d'Amazon API Gateway](#).

## sous-réseau

Plage d'adresses IP dans votre VPC. Un sous-réseau doit se trouver dans une seule zone de disponibilité.

## contrôle de supervision et acquisition de données (SCADA)

Dans le secteur manufacturier, un système qui utilise du matériel et des logiciels pour surveiller les actifs physiques et les opérations de production.

## chiffrement symétrique

Algorithme de chiffrement qui utilise la même clé pour chiffrer et déchiffrer les données.

## tests synthétiques

Tester un système de manière à simuler les interactions des utilisateurs afin de détecter les problèmes potentiels ou de surveiller les performances. Vous pouvez utiliser [Amazon CloudWatch Synthetics](#) pour créer ces tests.

## invite du système

Technique permettant de fournir un contexte, des instructions ou des directives à un [LLM](#) afin d'orienter son comportement. Les instructions du système aident à définir le contexte et à établir des règles pour les interactions avec les utilisateurs.

# T

## tags

Key-value des paires qui agissent comme des métadonnées pour organiser vos AWS ressources. Les balises peuvent vous aider à gérer, identifier, organiser, rechercher et filtrer des ressources. Pour plus d'informations, veuillez consulter la rubrique [Balisage de vos AWS ressources](#).

## variable cible

La valeur que vous essayez de prédire dans le cadre du ML supervisé. Elle est également qualifiée de variable de résultat. Par exemple, dans un environnement de fabrication, la variable cible peut être un défaut du produit.

## liste de tâches

Outil utilisé pour suivre les progrès dans un runbook. Liste de tâches qui contient une vue d'ensemble du runbook et une liste des tâches générales à effectuer. Pour chaque tâche générale, elle inclut le temps estimé nécessaire, le propriétaire et l'avancement.

## environnement de test

Voir [environnement](#).

## entraînement

Pour fournir des données à partir desquelles votre modèle de ML peut apprendre. Les données d'entraînement doivent contenir la bonne réponse. L'algorithme d'apprentissage identifie des modèles dans les données d'entraînement, qui mettent en correspondance les attributs des données d'entrée avec la cible (la réponse que vous souhaitez prédire). Il fournit un modèle de ML qui capture ces modèles. Vous pouvez alors utiliser le modèle de ML pour obtenir des prédictions sur de nouvelles données pour lesquelles vous ne connaissez pas la cible.

## outil

Fonction ou API qu'un [agent](#) peut invoquer pour effectuer des opérations dans des systèmes externes.

## passerelle de transit

Hub de transit de réseau que vous pouvez utiliser pour relier vos VPC et vos réseaux sur site. Pour plus d'informations, voir [Qu'est-ce qu'une passerelle de transit](#) dans la AWS Transit Gateway documentation.

## flux de travail basé sur jonction

Approche selon laquelle les développeurs génèrent et testent des fonctionnalités localement dans une branche de fonctionnalités, puis fusionnent ces modifications dans la branche principale. La branche principale est ensuite intégrée aux environnements de développement, de préproduction et de production, de manière séquentielle.

## accès sécurisé

Accorder des autorisations à un service que vous spécifiez pour effectuer des tâches au sein de votre organisation AWS Organizations et dans ses comptes en votre nom. Le service de confiance crée un rôle lié au service dans chaque compte, lorsque ce rôle est nécessaire, pour effectuer des tâches de gestion à votre place. Pour plus d'informations, consultez la section [Utilisation AWS Organizations avec d'autres AWS services](#) dans la AWS Organizations documentation.

## réglage

Pour modifier certains aspects de votre processus d'entraînement afin d'améliorer la précision du modèle de ML. Par exemple, vous pouvez entraîner le modèle de ML en générant un ensemble d'étiquetage, en ajoutant des étiquettes, puis en répétant ces étapes plusieurs fois avec différents paramètres pour optimiser le modèle.

## équipe de deux pizzas

Une petite DevOps équipe que vous pouvez nourrir avec deux pizzas. Une équipe de deux pizzas garantit les meilleures opportunités de collaboration possible dans le développement de logiciels.

# U

## incertitude

Un concept qui fait référence à des informations imprécises, incomplètes ou inconnues susceptibles de compromettre la fiabilité des modèles de ML prédictifs. Il existe deux types d'incertitude : l'incertitude épistémique est causée par des données limitées et incomplètes, alors que l'incertitude aléatoire est causée par le bruit et le caractère aléatoire inhérents aux données.

## tâches indifférenciées

Également connu sous le nom de « levage de charges lourdes », ce travail est nécessaire pour créer et exploiter une application, mais qui n'apporte pas de valeur directe à l'utilisateur final ni

d'avantage concurrentiel. Les exemples de tâches indifférenciées incluent l'approvisionnement, la maintenance et la planification des capacités.

environnements supérieurs

Voir [environnement](#).

## V

mise à vide

Opération de maintenance de base de données qui implique un nettoyage après des mises à jour incrémentielles afin de récupérer de l'espace de stockage et d'améliorer les performances.

contrôle de version

Processus et outils permettant de suivre les modifications, telles que les modifications apportées au code source dans un référentiel.

Appairage de VPC

Connexion entre deux VPC qui vous permet d'acheminer le trafic à l'aide d'adresses IP privées. Pour plus d'informations, veuillez consulter la rubrique [Qu'est-ce que l'appairage de VPC ?](#) dans la documentation Amazon VPC.

vulnérabilités

Défaut logiciel ou matériel qui compromet la sécurité du système.

## W

cache actif

Cache tampon qui contient les données actuelles et pertinentes fréquemment consultées. L'instance de base de données peut lire à partir du cache tampon, ce qui est plus rapide que la lecture à partir de la mémoire principale ou du disque.

données chaudes

Données rarement consultées. Lorsque vous interrogez ce type de données, des requêtes modérément lentes sont généralement acceptables.

## fonction de fenêtre

Fonction SQL qui effectue un calcul sur un groupe de lignes liées d'une manière ou d'une autre à l'enregistrement en cours. Les fonctions de fenêtre sont utiles pour traiter des tâches, telles que le calcul d'une moyenne mobile ou l'accès à la valeur des lignes en fonction de la position relative de la ligne en cours.

## charge de travail

Ensemble de ressources et de code qui fournit une valeur métier, par exemple une application destinée au client ou un processus de backend.

## flux de travail

Groupes fonctionnels d'un projet de migration chargés d'un ensemble de tâches spécifique. Chaque flux de travail est indépendant, mais prend en charge les autres flux de travail du projet. Par exemple, le flux de travail du portefeuille est chargé de prioriser les applications, de planifier les vagues et de collecter les métadonnées de migration. Le flux de travail du portefeuille fournit ces actifs au flux de travail de migration, qui migre ensuite les serveurs et les applications.

## VER

Voir [écrire une fois, lire plusieurs](#).

## WQF

Voir le [cadre AWS de qualification de la charge](#) de travail.

## écrire une fois, lire plusieurs (WORM)

Modèle de stockage qui écrit les données une seule fois et empêche leur suppression ou leur modification. Les utilisateurs autorisés peuvent lire les données autant de fois que nécessaire, mais ils ne peuvent pas les modifier. Cette infrastructure de stockage de données est considérée comme [immuable](#).

## Z

### exploit Zero-Day

Une attaque, généralement un logiciel malveillant, qui tire parti d'une [vulnérabilité de type « jour zéro »](#).

---

## vulnérabilité de type « jour zéro »

Une faille ou une vulnérabilité non atténuée dans un système de production. Les acteurs malveillants peuvent utiliser ce type de vulnérabilité pour attaquer le système. Les développeurs prennent souvent conscience de la vulnérabilité à la suite de l'attaque.

### invite Zero-Shot

Fournir à un [LLM](#) des instructions pour effectuer une tâche, mais aucun exemple (plans) pouvant aider à la guider. Le LLM doit utiliser ses connaissances pré-entraînées pour gérer la tâche. L'efficacité de l'invite zéro dépend de la complexité de la tâche et de la qualité de l'invite. Voir également les instructions [en quelques clics](#).

### application zombie

Application dont l'utilisation moyenne du processeur et de la mémoire est inférieure à 5 %. Dans un projet de migration, il est courant de retirer ces applications.

---

Les traductions sont fournies par des outils de traduction automatique. En cas de conflit entre le contenu d'une traduction et celui de la version originale en anglais, la version anglaise prévaudra.