

Ingenierie des donnees sur Google Cloud

Prix : 3 160 € HT

Durée : 4 jours

Code de Référence : GCP200DE

Catalogue Google Cloud Platform

Cette formation cours utilise des presentations, des demonstrations et des travaux pratiques pour vous montrer comment concevoir des systemes de traitement de donnees, creer des pipelines de donnees de bout en bout, analyser des donnees et implementer le machine learning.

Objectifs de la formation

Cette formation Google Cloud permet aux participants d'acquérir les compétences suivantes:

- Concevoir des systemes de traitement de donnees evolutifs dans Google Cloud
- Differencier les architectures de donnees et implementer les concepts de lakehouse et de pipelines de donnees
- Construire et gerer des pipelines de donnees robustes en streaming et en batch
- Utiliser les outils IA/ML pour optimiser les performances et obtenir des informations sur les processus et les donnees

Public

Cette formation Google Cloud s'adresse aux ingenieurs de donnees, analystes de donnees, architectes de données.

Cette formation est accessible aux personnes en situation de handicap, contactez-nous pour en savoir plus.

Prérequis

Pour tirer le meilleur parti de ce cours, les participants doivent disposer des éléments suivants :

- Compréhension des principes d'ingenierie des donnees, y compris les processus ETL/ELT, la modelisation des donnees et les formats de donnees courants (Avro, Parquet, JSON)
- Familiarite avec les concepts d'architecture de donnees, en particulier les entrepots de donnees (Data Warehouses) et les lacs de donnees (Data Lakes)
- Maitrise de SQL pour l'interrogation des donnees
- Maitrise d'un langage de programmation courant (Python recommande)
- Familiarite avec l'utilisation des interfaces de ligne de commande (CLI)

- Familiarité avec les concepts et services de base de Google Cloud (Compute, Storage et gestion des identités)

Vous souhaitez faire vérifier vos prérequis ? Contactez-nous pour l'organisation d'un entretien téléphonique avec un de nos consultants formateurs.

Pour une efficacité renforcée, le nombre de participants est limitée à 12. Le maintien des sessions est conditionné à un minimum de 3 participants.

Programme de la formation

Module 1 : Taches et composants de l'ingénierie des données

- Expliquer le rôle d'un ingénieur de données
- Comprendre les différences entre une source de données et un récepteur de données
- Expliquer les différents types de formats de données
- Expliquer les options de solutions de stockage sur Google Cloud
- Découvrir les options de gestion des métadonnées sur Google Cloud
- Comprendre comment partager facilement des jeux de données avec Analytics Hub
- Comprendre comment charger des données dans BigQuery à l'aide de la console Google Cloud ou de la CLI gcloud

Lab : Chargement de données dans BigQuery

Quiz

Module 2 : Réplication et migration de données

- Expliquer l'architecture de base de réplication et de migration de données de Google Cloud
- Comprendre les options et les cas d'utilisation de l'outil de ligne de commande gcloud
- Expliquer la fonctionnalité et les cas d'utilisation de Storage Transfer Service
- Expliquer la fonctionnalité et les cas d'utilisation de Transfer Appliance
- Comprendre les fonctionnalités et le déploiement de Datastream

Module 3 : Le modèle de pipeline de données d'extraction et de chargement

- Expliquer le schéma d'architecture de base d'extraction et de chargement
- Comprendre les options de l'outil de ligne de commande bq
- Expliquer la fonctionnalité et les cas d'utilisation du service de transfert de données BigQuery
- Expliquer la fonctionnalité et les cas d'utilisation de BigLake en tant que modèle sans extraction-chargement

Lab : BigLake : Démarrage rapide

Quiz

Module 4 : Le modèle de pipeline de données d'extraction, de chargement et de transformation

- Expliquer le schéma d'architecture de base d'extraction, de chargement et de transformation
- Comprendre un pipeline ELT courant sur Google Cloud
- Découvrir les capacités de scripting SQL et de planification de BigQuery
- Expliquer la fonctionnalité et les cas d'utilisation de Dataform

Lab : Créer et exécuter un workflow SQL dans Dataform

Quiz

Module 5: Le modèle de pipeline de données d'extraction, de transformation et de chargement

- Expliquer le schéma d'architecture de base d'extraction, de transformation et de chargement
- Découvrir les outils d'interface graphique sur Google Cloud utilisés pour les pipelines de données ETL
- Expliquer le traitement des données en batch avec Dataproc
- Apprendre à utiliser Dataproc Serverless pour Spark pour l'ETL
- Expliquer les options de traitement des données en streaming
- Expliquer le rôle que joue Bigtable dans les pipelines de données

Lab : Utiliser Dataproc Serverless pour Spark pour charger BigQuery (optionnel)

Lab : Créer un pipeline de données en streaming pour un tableau de bord en temps réel avec Dataflow

Quiz

Module 6 : Techniques d'automatisation

- Expliquer les modèles d'automatisation et les options disponibles pour les pipelines
- Découvrir Cloud Scheduler et Workflows
- Découvrir Cloud Composer
- Découvrir Cloud Run Functions
- Expliquer la fonctionnalité et les cas d'utilisation d'automatisation pour Eventarc

Lab : Utiliser Cloud Run Functions pour charger BigQuery (optionnel)

Quiz

Module 7 : Introduction à l'ingénierie des données moderne sur Google Cloud

- Comparer et contraster les architectures de lac de données, d'entrepôt de données et de lakehouse de données
- Évaluer les avantages de l'approche lakehouse

Quiz

Module 8 : Construire un lakehouse de données avec Cloud Storage, les formats ouverts et BigQuery

- Discuter des options de stockage de données, y compris Cloud Storage pour les fichiers, les formats de table ouverts comme Apache Iceberg, BigQuery pour les données analytiques et AlloyDB pour les données opérationnelles
- Comprendre le rôle d'AlloyDB pour les cas d'utilisation de données opérationnelles

Quiz

Lab : Requête fédérée avec BigQuery

Module 9 : Moderniser les entrepôts de données avec BigQuery et BigLake

- Expliquer pourquoi BigQuery est une solution d'entreposage de données évolutive sur Google Cloud

- Discuter des concepts de base de BigQuery
- Comprendre le rôle de BigLake dans la création d'une architecture lakehouse unifiée et son intégration avec BigQuery pour les données externes
- Apprendre comment BigQuery interagit nativement avec les tables Apache Iceberg via BigLake

Quiz

Lab : Interroger des données externes et des tables Iceberg

Module 10 : Modèles avancés de lakehouse et gouvernance des données

- Implémenter des pratiques robustes de gouvernance et de sécurité des données sur la plateforme de données unifiée, y compris la protection des données sensibles et la gestion des métadonnées
- Explorer l'analytique avancée et le machine learning directement sur les données du lakehouse

Quiz

Module 11 : Labs et bonnes pratiques

- Renforcer les principes fondamentaux de la plateforme de données de Google Cloud

Lab : Démarrer avec BigQuery ML

Lab : Recherche vectorielle avec BigQuery

Module 12 : Quand choisir les pipelines de données en batch

- Expliquer le rôle critique d'un ingénieur de données dans le développement et la maintenance des pipelines de données en batch
- Décrire les composants de base et le cycle de vie typique des pipelines de données en batch, de l'ingestion à la consommation en aval
- Analyser les défis courants du traitement de données en batch, tels que le volume de données, la qualité, la complexité et la fiabilité, et identifier les services Google Cloud clés qui peuvent les résoudre

Quiz

Module 13 : Concevoir et construire des pipelines de données en batch évolutifs

- Concevoir des pipelines de données en batch évolutifs pour l'ingestion et la transformation de données à haut volume
- Optimiser les jobs en batch pour un haut débit et une efficacité des coûts en utilisant diverses techniques de gestion des ressources et d'ajustement des performances

Quiz

Lab : Construire un pipeline de données en batch simple avec Serverless pour Apache Spark (optionnel)

Lab : Construire un pipeline de données en batch simple avec l'interface Dataflow Job Builder (optionnel)

Module 14 : Contrôler la qualité des données dans les pipelines de données en batch

- Développer des règles de validation des données et une logique de nettoyage pour assurer la qualité des données dans les pipelines en batch
- Implémenter des stratégies pour gérer l'évolution des schémas et effectuer la déduplication des données dans les grands jeux de données

Lab : Valider la qualité des données dans un pipeline en batch avec Serverless pour Apache Spark (optionnel)

Quiz

Module 15 : Orchestrer et surveiller les pipelines de données en batch

- Orchestrer des workflows de pipelines de données en batch complexes pour une planification efficace et un suivi de lignage
- Implémenter une gestion robuste des erreurs, une surveillance et une observabilité pour les pipelines de données en batch

Lab : Construire des pipelines en batch dans Cloud Data Fusion

Quiz

Module 16 : Construire des pipelines de données en streaming sur Google Cloud

- Introduire les objectifs d'apprentissage du cours et le scénario qui sera utilisé pour apporter un apprentissage pratique à la construction de pipelines de données en streaming
- Décrire le concept de pipelines de données en streaming, les défis associés et le rôle de ces pipelines dans le processus d'ingénierie des données

Lab : Utiliser des modèles d'IA sur Kubeflow

Module 17 : Cas d'utilisation du streaming et architectures de référence

- Comprendre les différents cas d'utilisation du streaming et leurs applications, y compris le Streaming ETL, le Streaming IA/ML, les applications de streaming et le Reverse ETL
- Identifier et décrire les architectures types courantes pour les données en streaming, y compris le Streaming ETL, le Streaming IA/ML, les applications de streaming et le Reverse ETL

Quiz

Module 18 : Plongée approfondie dans les produits

- Pub/Sub et Managed Service for Apache Kafka : Définir les concepts de messagerie, savoir quand utiliser Pub/Sub ou Managed Service for Apache Kafka
- Dataflow : Décrire le service et les défis avec les données en streaming, construire et déployer un pipeline de streaming
- BigQuery : Explorer les différentes méthodes d'ingestion de données, utiliser les requêtes continues BigQuery, BigQuery ETL et le reverse ETL, configurer le streaming Pub/Sub vers BigQuery, architecturer les pipelines de streaming BigQuery
- Bigtable : Décrire la vue d'ensemble du mouvement et de l'interaction des données, établir un pipeline de streaming de Dataflow vers Bigtable, analyser le flux de données continu Bigtable pour les tendances avec BigQuery, synchroniser l'analyse des tendances dans l'application utilisateur

Lab : Streamer des données avec des pipelines – Cas d'utilisation Esports (optionnel)

Quiz

Lab : Utiliser Apache Beam et Bigtable pour enrichir les données de contenu téléchargeables (DLC) esports

Quiz

Lab : Streamer des données e-sports avec Pub/Sub et BigQuery

Quiz

Lab : Surveiller le chat e-sports avec Streamlit

Quiz

Méthodes pédagogiques

Des exercices pratiques et des démonstrations vous permettront de mettre en pratique les notions théoriques présentées.

Méthodes d'évaluation des acquis

Afin d'évaluer l'acquisition de vos connaissances et compétences, il vous sera envoyé un formulaire d'auto-évaluation, qui sera à compléter en amont et à l'issue de la formation.

Un certificat de réalisation de fin de formation est remis au stagiaire lui permettant de faire valoir le suivi de la formation.