



COLLÈGE
ROSEMONT

Big Data : Systèmes de gestion de données

Par : Abderrazak Sahraoui

Sommaire

- Base de données
- Entrepôt de données
- Lac de données
- Entrepôt lac de données

- Quand utiliser une base de donnée vs un entrepôt de données vs lac de données?
- Quel est le rôle de l'architecte dans la construction des systèmes de données?
- Quelle est la différence entre un schéma de base de données et schéma d'un entrepôt?
- Une donnée voyage-t-elle durant son cycle de vie entre BD, ED et LD ?
- Comment déterminer la fraîcheur d'une donnée ?
- Quelle est la source de données pour chaque structure ?

Base de données



https://www.youtube.com/watch?v=-bSkREem8dM&ab_channel=AlexTheAnalyst

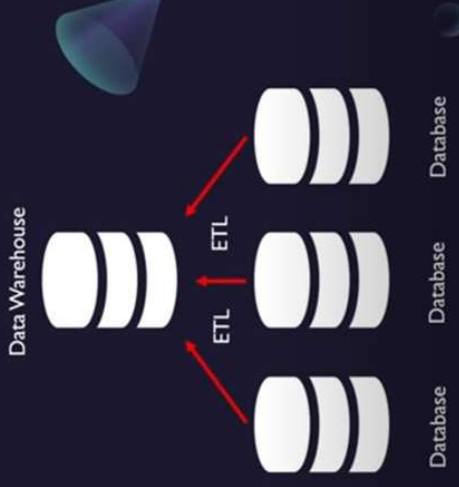
- Une base de données permet d'enregistrer des données provenant d'un processus OLTP (Online Transaction Process). Autrement dit, elle permet de capturer des données fraîchement créés par une application.
- Les données sont stockées dans les bases de données relationnelles sous forme de table de plusieurs colonnes et plusieurs lignes.
- Le schéma relationnel d'une table peut être aisément modifié comparé à s structures. Il est possible d'ajouter de nouvelles colonnes à une table ou d'en supprimer ou modifier.

Entrepôt de données

What is a Data Warehouse?

RELATIONAL DATABASE

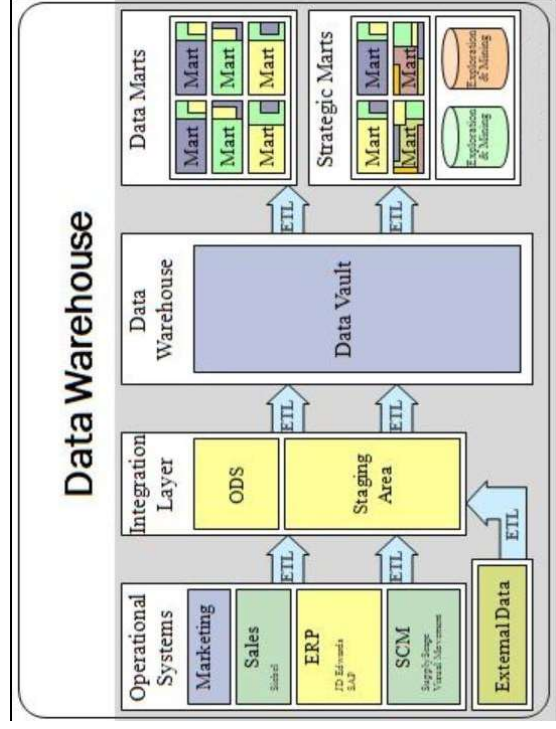
- Designed for analytical processing (OLAP)
- Data is refreshed from source systems – stores current and historical
- Data is summarized
- Rigid Schema (how the data is organized)



- Un entrepôt de données permet d'entreposer des données provenant de différentes sources de données. Lesquelles données sont destinées à être analysées par un processus OLAP (Online Analytics Process)
- Le but d'un entrepôt de données est de fournir une référence unique pour un ensemble de données pouvant servir dans la prise de décisions au sein de l'entreprise, et d'offrir les outils nécessaires aux processus analytiques BI (Business intelligence ou Informatique décisionnelle).

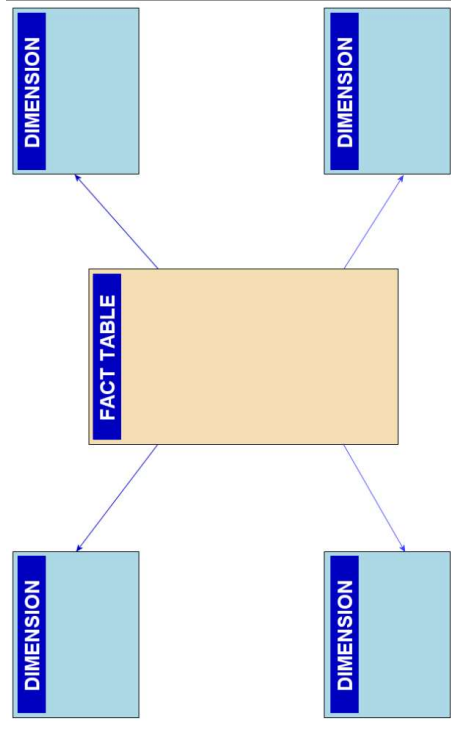
Entrepôt de données

- En amont, Les données arrivent à l'entrepôt par le biais d'un processus ETL (Extract, Transforme et Load). Les données sont extraites de sources localisées dans des systèmes transactionnels en production. Les données sont épurées ou transformées par filtrage, codage et certification.
- Les données de l'entrepôt peuvent être conservées sous deux formes :
 - sous forme élémentaire et détaillée.
 - sous forme agrégée selon des axes ou des dimensions d'analyse prévues. il est impossible de retrouver le détail et la profondeur des indicateurs une fois ceux-ci agrégés. (par exemple, si l'on a agrégé les résultats par mois, il ne sera plus possible de faire une analyse par journée).
- En aval, les données peuvent être restituées aux usagers par des outils OLAP de :
 - requêtes ou reporting,
 - cubes ou hypercubes,
 - fouille de données.



https://fr.wikipedia.org/wiki/Entrep%C3%B4t_de_donn%C3%A9es

Entrepôt de données



- Les entrepôts utilisent le modèle dit en étoile ou les tables sont réparties en deux catégories : tables de faits et tables de dimension.
- Chaque modèle en étoile est constitué d'une table centrale de faits contenant les mesures comme montant, quantité, etc. et de plusieurs tables de dimension comme le temps (jour, mois, année) nomenclature (famille de produit, sous-famille, ...) segmentation clientèle (sexe, tranche âge,...)
- La jointure dans un modèle en étoile entre table de faits et tables de dimension est facilitée (optimisée) par la présence d'une clé calculée à partir des clés des tables de dimension ce qui facilite l'analytique.
- Le modèle dit en flocon est une variante du modèle en étoile où les tables de dimensions sont normalisées évitant ainsi le redondance et permettant un gain d'espace de l'ordre de 5 à 10%.

[https://fr.wikipedia.org/wiki/%C3%89toile_\(mod%C3%A8le_de_dimension%C3%A9es\)](https://fr.wikipedia.org/wiki/%C3%89toile_(mod%C3%A8le_de_dimension%C3%A9es))

Comparatif

Key Differences

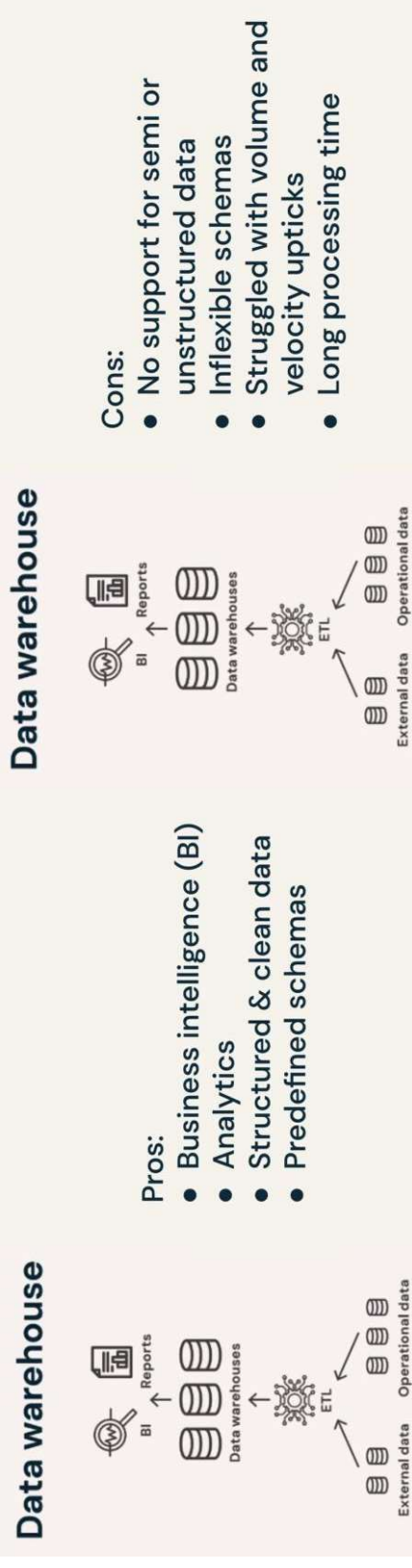
- Databases are designed for Transactions, Data Warehouses are designed for analytics and reporting
- Databases data is fresh and detailed, Data Warehouses data is refreshed periodically and is summarized
- Databases work slowly for querying large amounts of data and can slow down transactional processes, Data Warehouses don't interfere with any processes and are generally faster



Comparatif

Caractéristique	Base de données de production	Data warehouses	Datamarts
Opération	gestion courante, production	référentiel, analyse ponctuelle	analyse récurrente, outil de pilotage, support à la décision
Modèle de données	entité/relation	3NF, étoile, flocon	étoile, flocon
Normalisation	fréquente	maximum	rare (redondance d'information)
Données	actuelles, brutes, détaillées	historisées, détaillées	historisées, agrégées
Mise à jour	immédiate, temps réel	souvent différée, périodique	souvent différée, périodique
Niveau de consolidation	faible	faible	élevé
Perception	verticale	transverse	horizontale
Opérations	lectures, insertions, mises à jour, suppressions	lectures, insertions, mises à jour	lectures, insertions, mises à jour, suppressions
Taille	en gigaoctets	en téraoctets	en gigaoctets

Data Warehouse



Lac de données

What is a Data Lake?

RELATIONAL DATABASE

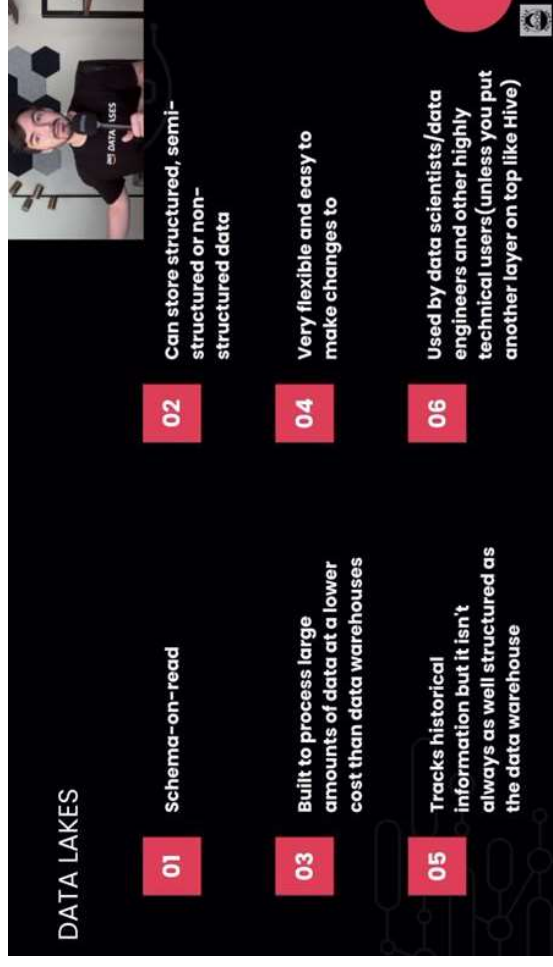
- Designed to capture raw data (structured, semi-structured, unstructured)
- Made for large amounts of data
- Used for ML and AI in its current state or for Analytics with processing
- Can organize and put into Databases or Data Warehouses



https://fr.wikipedia.org/wiki/Lac_de_donn%C3%A9es

- Un lac de donnée permet le stockage rapide de données massives hétérogènes dans leur format original ou avec peu de transformation.
 - Données structurées issues de bd relationnelles.
 - Données issues de bases NoSQL
 - Données semi-structurées (fichiers CSV, journaux, xml, json,...)
 - Données non structurées (emails, documents, pdf
 - Fichiers blob (images, audio, vidéo)
- Les lacs sont utilisés par des ingénieurs de données et des scientifiques de données pour des applications en apprentissage machine et intelligence artificielle.
- Lorsqu'une donnée arrive au lac, elle se verra dotée d'un identifiant et de balises de métadonnées. Lorsqu'un besoin se présente, le Data Lake est parcouru pour y rechercher des informations pertinentes. L'analyse de ces données permet alors d'apporter de la valeur et de répondre à ce besoin.
- Le stockage se fait en utilisant l'architecture d'un cluster Hadoop.
- Les données peuvent être conservées dans le lac pour un usage ultérieur non prédéterminé.

Lac de données



https://www.youtube.com/watch?v=ExpRL0m9BcA&ab_channel=SeattleDataGuy

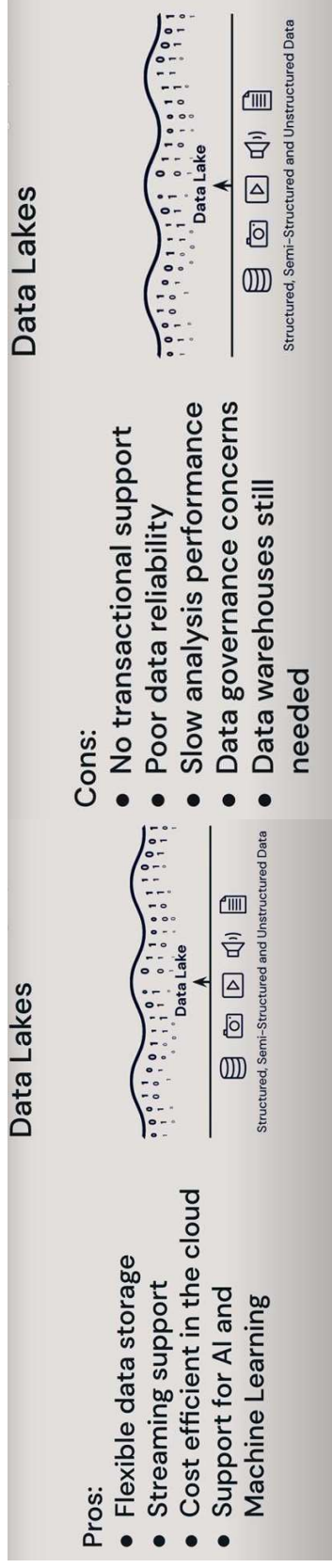
- Les avantages des lacs de données sont :
 - la rationalisation du stockage des données,
 - la réduction des coûts de stockage,
 - et la facilitation de l'accès pour l'analyse et la prise de décisions d'une façon globale.
- Les inconvénients sont :
 - la difficulté à conserver un lac de données propre et organisé,
 - la difficulté à organiser et maintenir une gouvernance des données efficace,
 - le temps nécessaire à traiter et analyser les données stockées à l'état brut.
 - L'expertise requise pour rechercher, analyser et traiter les données de manière pertinente et créatrice de valeur, souvent confiées aux Data Scientists
 - la sécurité, la confidentialité et les problématiques liées aux données personnelles et au respect des réglementations.
- Plusieurs environnements fournissent des services complets permettant la gestion d'un lac de données. La plupart d'entre eux sont basés sur la technologie Hadoop et fournissent des installations en local (MapR, Cloudera, Hortonworks) ou dans le Cloud (Microsoft Azure, Google Cloud Platform, Amazon S3)

Data warehouse vs. data lake

•

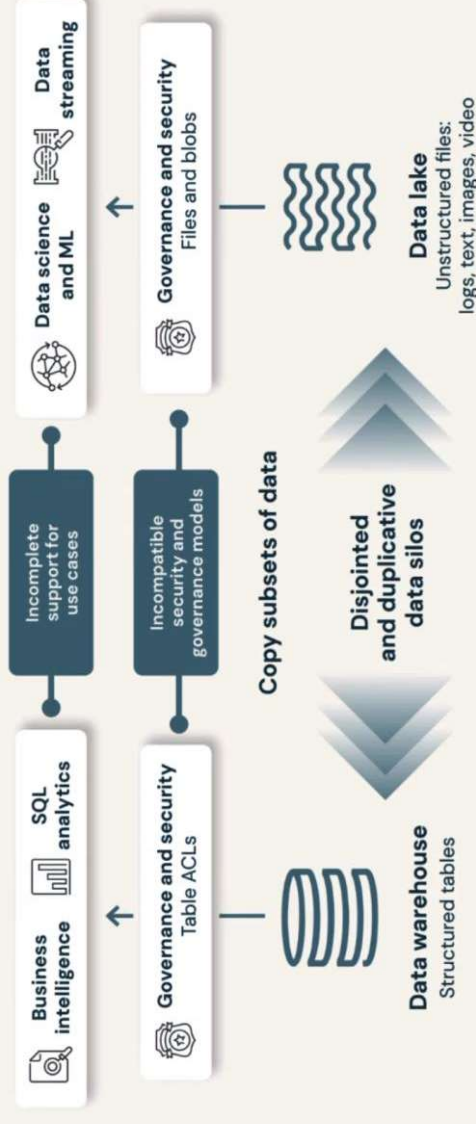
	DATA WAREHOUSE	DATA LAKE
DATA TYPES	Structured, processed data from operational databases, applications and transactional systems	Structured, semistructured and unstructured data from sensors, apps, websites, etc.
PURPOSE	Predefined purpose for business intelligence, batch reporting and data visualization	May not have a predefined purpose; typically used for machine learning, deep analysis and discovery
USERS	Data engineers, business analysts, data analysts	Data engineers, data scientists
SCHEMA POSITION	Schema-on-write	Schema-on-read
BENEFITS	Categorized historical data stored in a single repository with ease of access for the end user	Data stored in its native format, allowing flexibility for data scientists to analyze and develop models from diverse data sources

Lac de données



Entrepôt et lac de données

Business required two disparate, incompatible data platforms



- Besoin d'une plateforme de données qui combine les avantages des deux solutions et élimine leurs inconvénients.

Databricks Lakehouse

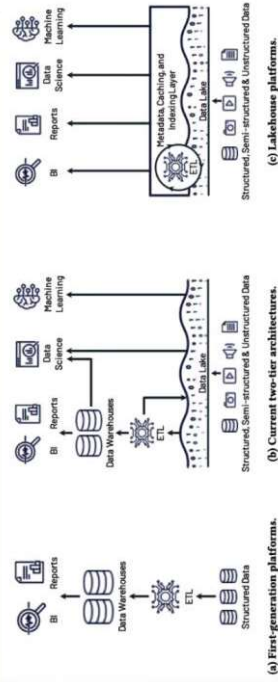


Figure 1: Evolution of data platform architectures to today's two-tier model (a-b) and the new Lakehouse model (c).

Lakehouse: A New Generation of Open Platforms that Unify Data Warehousing and Advanced Analytics. M. Armbrust, A. Ghodsi, R. Xin, M. Zaharia. 11th Annual Conference on Innovative Data Systems Research (CIDR '21), January 11–15, 2021, Online.

Lakehouse: A New Generation of Open Platforms that Unify Data Warehousing and Advanced Analytics

Michael Armbrust¹, Ali Ghodsi^{1,2}, Reynold Xin¹, Matei Zaharia^{1,2}

¹Databricks, NYC, Berkeley

²Stanford University

Abstract

This paper argues that the data warehouse architecture as we know it today will soon be replaced by a new architecture. The new architecture is based on open data lake data formats, such as Apache Parquet, (i) have first-class support for analytics, (ii) have first-class support for data science, and (iii) have first-class support for data engineering. Lakehouses can help address several major challenges with data warehouses, including data silos, data quality, and data security. We discuss how the industry is already moving toward lakehouses, and we argue that lakehouses will become the new standard for data storage and analytics. We also report on our experience with a lakehouse system using Parquet that is competitive with popular cloud data warehouses on TPC-DS.

1 Introduction

This paper argues that the data warehouse architecture as we know it today will soon be replaced by a new architecture. The new architecture is based on open data lake data formats, such as Apache Parquet, (i) have first-class support for analytics, (ii) have first-class support for data science, and (iii) have first-class support for data engineering. Lakehouses can help address several major challenges with data warehouses, including data silos, data quality, and data security. We discuss how the industry is already moving toward lakehouses, and we argue that lakehouses will become the new standard for data storage and analytics. We also report on our experience with a lakehouse system using Parquet that is competitive with popular cloud data warehouses on TPC-DS.

The history of data warehousing started with helping business decision support and business intelligence (BI). Data in these systems was typically stored in a single table, and the data was not structured. The data was then loaded into a data warehouse, which was a centralized system for storing and analyzing data. The data warehouse architecture was designed to support a wide range of analytics, from simple queries to complex data science workloads. However, the data warehouse architecture has several limitations, including data silos, data quality, and data security. Lakehouses can help address these limitations by providing a unified platform for storing and analyzing data.

A decade ago, the first generation system started a few novel challenges. First, they typically coupled compute and storage into a single layer, which made it difficult to scale compute independently of storage. Second, they only used a single format for storing data, which made it difficult to support a wide range of analytics. Third, they only supported a limited range of data types, which made it difficult to support a wide range of data science workloads. Lakehouses can help address these challenges by providing a unified platform for storing and analyzing data.

The second generation data analytics platforms started addressing all the new data lake data formats, but they still had several limitations. First, they only supported a limited range of data types, which made it difficult to support a wide range of data science workloads. Second, they only supported a limited range of analytics, which made it difficult to support a wide range of business intelligence workloads. Third, they only supported a limited range of data sources, which made it difficult to support a wide range of data integration scenarios. Lakehouses can help address these limitations by providing a unified platform for storing and analyzing data.

This work is published under the Creative Commons Attribution License (CC BY). Copyright 2021 Databricks. All rights reserved.

quality and governance downstream. In this architecture, a small subset of data in the lake would later be ETLED to a downstream data warehouse. The use of open formats also made it easier to integrate new data sources and analytics. From 2015 onwards, cloud data lakes such as S3, ADLS and GCS, have become the primary data storage for many enterprises. These data lakes have several advantages over traditional data warehouses, including lower cost, greater flexibility, and easier integration with other data sources. However, they also have several limitations, including data silos, data quality, and data security. Lakehouses can help address these limitations by providing a unified platform for storing and analyzing data.

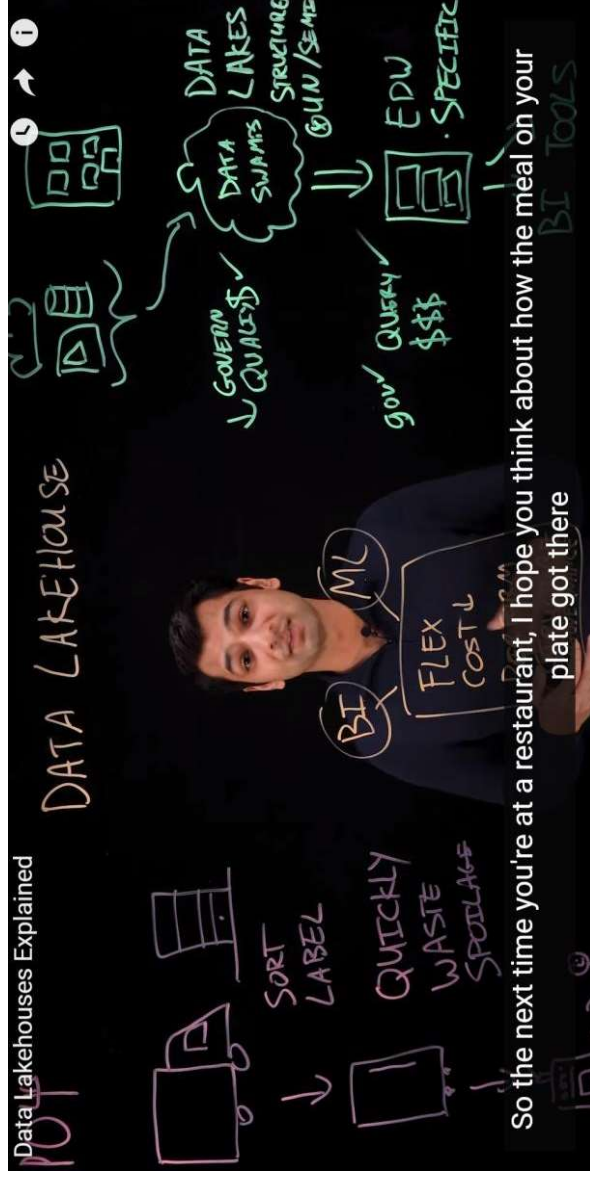
While the cloud data lake and warehouse architectures are naturally complementary, they are not mutually exclusive. In fact, many enterprises are using both architectures to support their data needs. The cloud data lake architecture is used for storing and analyzing large volumes of data, while the cloud data warehouse architecture is used for storing and analyzing smaller volumes of data. Lakehouses can help address the limitations of both architectures by providing a unified platform for storing and analyzing data.

Today's data architectures commonly suffer from four problems: **Availability.** Keeping the data lake and warehouse separate in the first place makes it difficult to ensure high availability. **Performance.** The data lake and warehouse architectures are not designed to support a wide range of analytics, which makes it difficult to achieve high performance. **Cost.** The data lake and warehouse architectures are not designed to support a wide range of data types, which makes it difficult to achieve low cost. **Security.** The data lake and warehouse architectures are not designed to support a wide range of data sources, which makes it difficult to achieve high security.

This is a step back compared to the first generation of analytics platforms, which were designed to support a wide range of analytics. The data lake and warehouse architectures are not designed to support a wide range of analytics, which makes it difficult to achieve high performance. Lakehouses can help address these limitations by providing a unified platform for storing and analyzing data.

Limited support for advanced analytics. Businesses want to use their data for a wide range of analytics, from simple queries to complex data science workloads. However, the data lake and warehouse architectures are not designed to support a wide range of analytics, which makes it difficult to achieve high performance. Lakehouses can help address these limitations by providing a unified platform for storing and analyzing data.

Lakehouse : Entrepôt lac de données



....

https://www.youtube.com/watch?v=Enu-EH7RHHM&ab_channel=IBMTTechnology

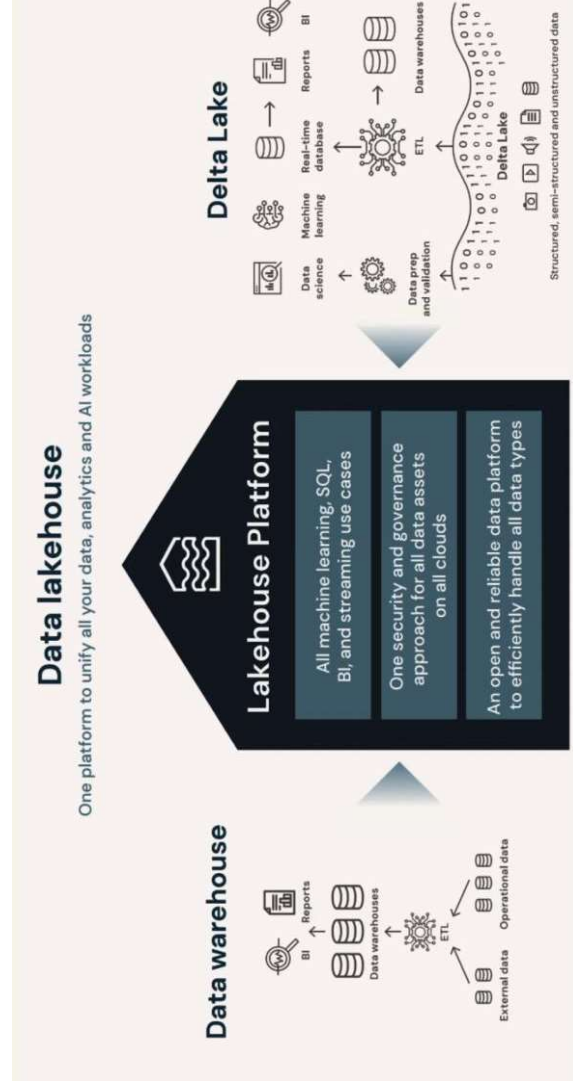
Databricks Lakehouse



https://www.youtube.com/watch?v=Enu-EH7RHHM&ab_channel=IBMTTechnology

Entrepôt lac de données

- Databricks offre une architecture hybride unifiant une plateforme d'entrepôt de données et une plateforme de lac de données.



Entrepôt lac de données

Key features of a data lakehouse.

- Transaction support
 - Schema enforcement and governance
 - Data governance
 - BI Support
 - Decoupled storage from compute
- Open storage formats
 - Support for diverse data types
 - Support for diverse workloads
 - End-to-end streaming