

Big Data : technologies

BIG DATA





Semifir

Les Technologies du Big Data

Ensembles des techniques et strategies permettant de répondre aux contraintes des grands **V** du Big Data.

Parradigmes du Big Data

Pour repondre aux besoins du Big Data, des sytemes permettant de depasser les limites des systemes traditionnels ont été developpés.

- Architectures scallables
- Solutions de stockages pour les données non structurées

- Architectures distribuables et massivement parallèles
 - Algorithmes distribués et parallèles
-

Scalabilité

Adapter la taille et/ou la puissance d'un système informatique pour répondre aux changements de la charge de travail.

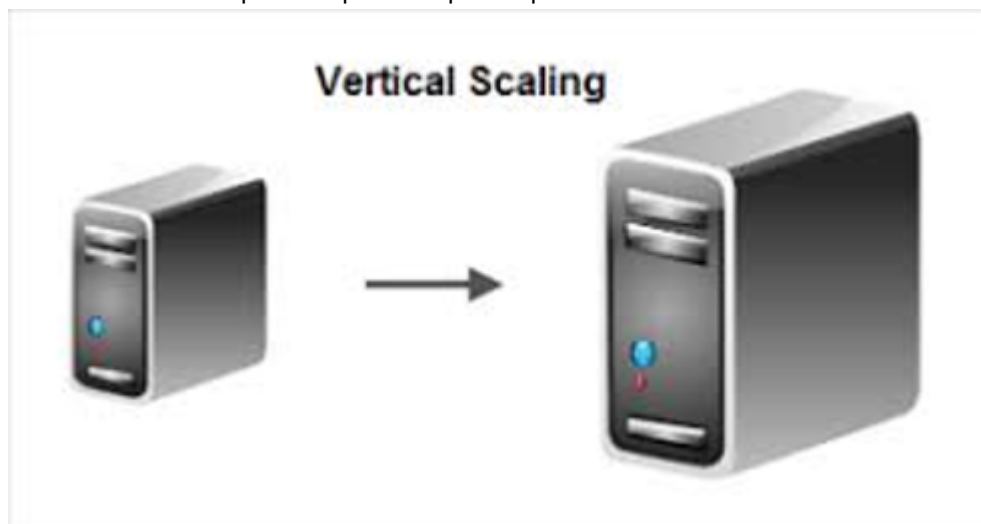
Il existe deux types de scalabilité :

- **Verticale** : Augmentation des ressources internes.
 - **Horizontale** : Augmentation des ressources externes.
-

Scalabilité verticale

Augmentation de la puissance (processeur, RAM, stockage) d'un système informatique.

- C'est la solution la plus simple et la plus rapide à mettre en oeuvre.

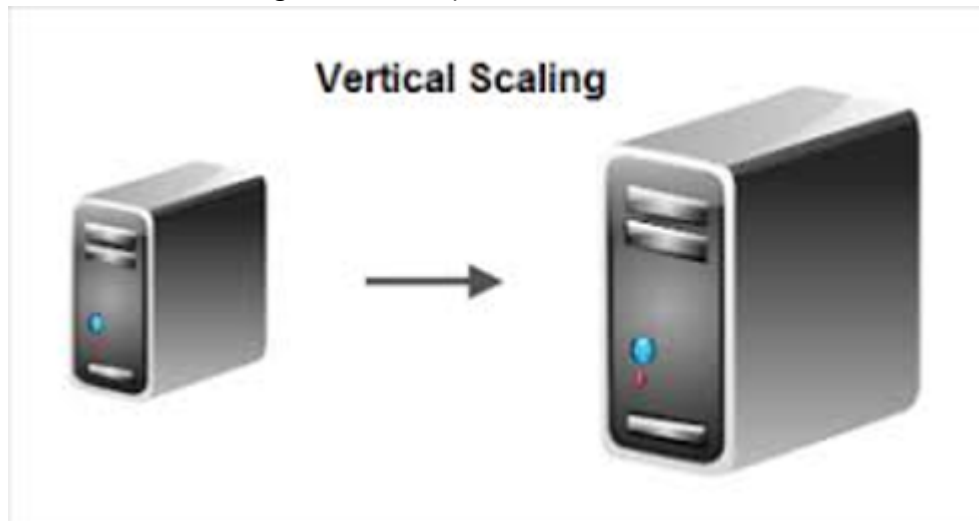


- Fréquemment utilisée dans les systèmes traditionnels.
-

Limites de la scalabilité verticale

L'augmentation de la puissance d'un système informatique est limitée par :

- Cout du materiels (augmentation exponentielle)



- Faible adaptabilité aux changements de la charge de travail
- Problemes en cas de pannes (Single point of failure)

Scalabilité horizontale

Augmentation du nombre de machine de faible puissance pour augmenter la puissance globale.

- C'est la solution économique la plus adaptée.



- Pseudo-linéarité des performances.
- Decoupages et repliations des données.

Limites de la scalabilité horizontale

L'augmentation du nombre de machine de faible puissance est limitée par :

- Augmentation exponentielle des échanges



- Architectures reseaux complexes (cluster)
- Synchronisation des données

Modèles de stockage

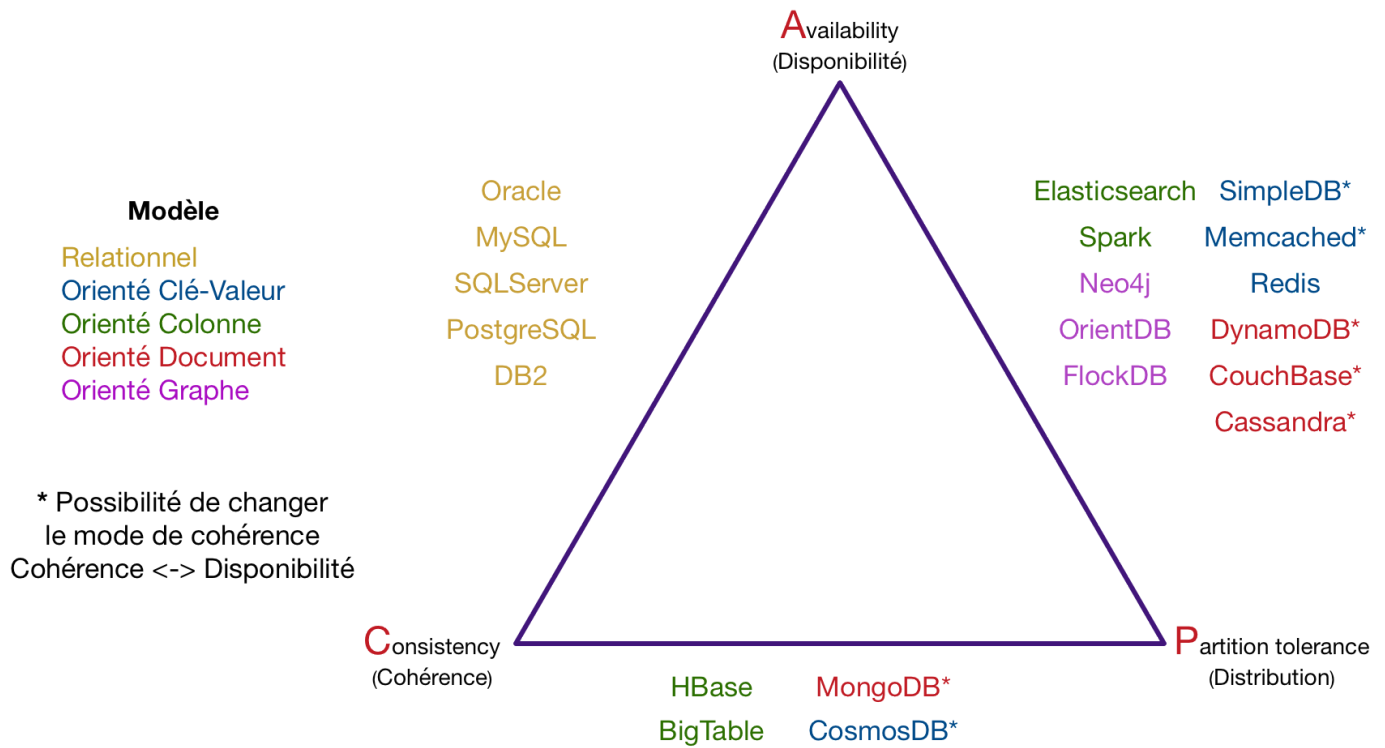
L'avènement du big data a nécessité l'adaptation des systèmes de stockage pour s'adapter à la quantités et aux types de données.

- Theoreme de CAP et big data
- Modeles NoSQL

Theoreme de CAP

Tous systemes ne peut garantir que 2 des 3 propriétés suivantes :

- **C**onsistency (Cohérence) : accès aux mêmes données à tous moments.
 - **A**vailability (Disponibilité) : accès à la lecture et l'écriture à tous moments.
 - **P**artition Tolerance (Tolérance aux pannes)
-



Stockage traditionnel

Les systèmes de stockage traditionnels utilisent des bases de données relationnelles.

- Données structurées et formatées



- Requêtes et manipulations simples
- Transaction qui respectent les principes ACID
- Large panel de solutions (longévité des technologies)

Limites des bases de données relationnelles

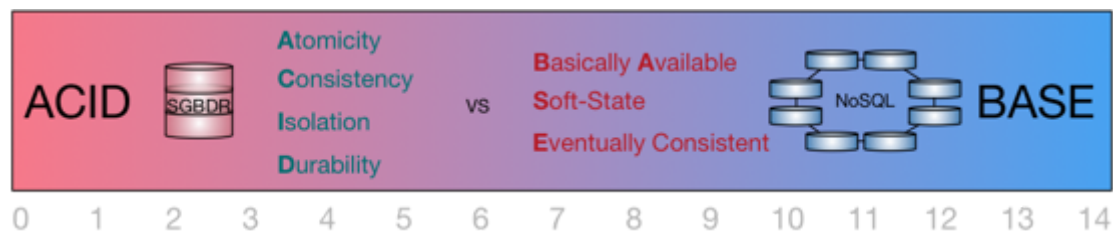
- Manque d'adaptabilité des schémas de données

- Limitées aux données structurées (<20%)
- Dépenses massives en temps et ressources en cas de modifications de la structure des données.

Avenement du NoSQL

- Développer pour pallier aux limitations des bases de données relationnelles.
- Prise en charge des données non structurées.
- Adaptabilité des schemas de données.
- Abandon des principes ACID au profit des principes BASE.

Les principes BASE

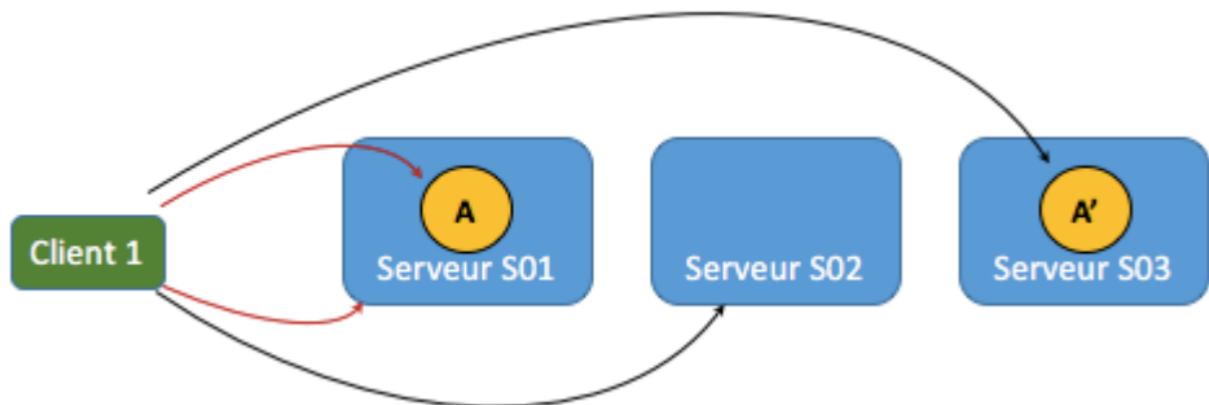


- **B**asically **A**vailable
- **S**oft-State
- **E**ventual Consistency

Basically Available

Notion de disponibilité des données à tout moment.

- Le système doit pouvoir répondre aux requêtes de tout utilisateur même en cas de pannes.



- Les requêtes peuvent être obsolètes. (pert Isolation ACID).

Soft-State

Notion que les données sont dans un flux d'utilisation constante.

- Les données peuvent être utilisées par plusieurs utilisateurs en même temps.

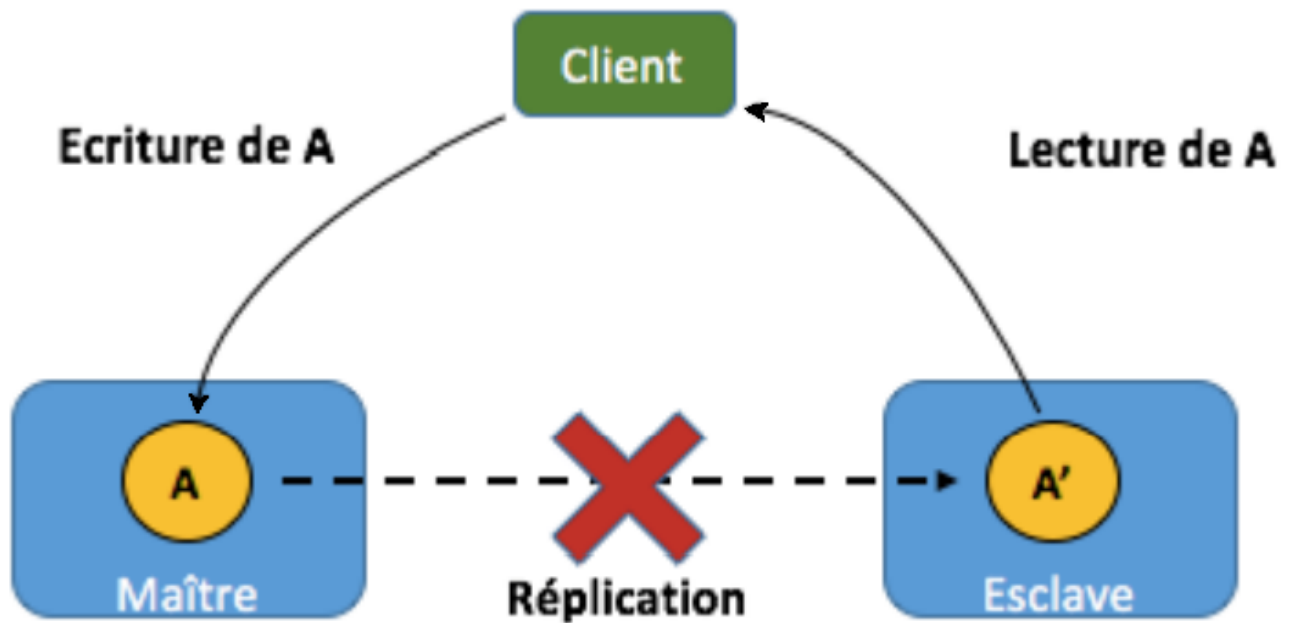


- La cohérence n'est pas garantie (pert Cohérence ACID).

Eventual Consistency

Notion qu'au bout d'un certain temps le système sera cohérent.

- La synchronisation des données est faite en arrière plan.



- Les données peuvent être obsolètes (pert Cohérence ACID).

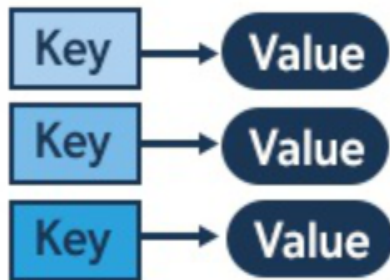
Les modèles NoSQL

A partir de ces principes, quatre modèles de stockage NoSql ont été développés :

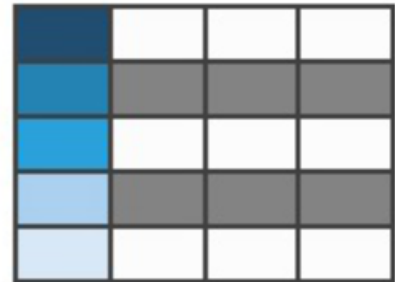
- Clé/valeur

NoSQL

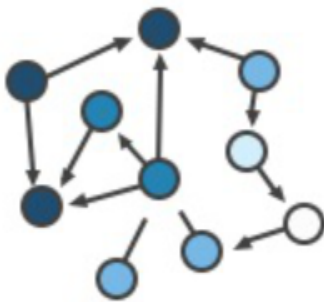
Key-Value



Column-Family



Graph



Document

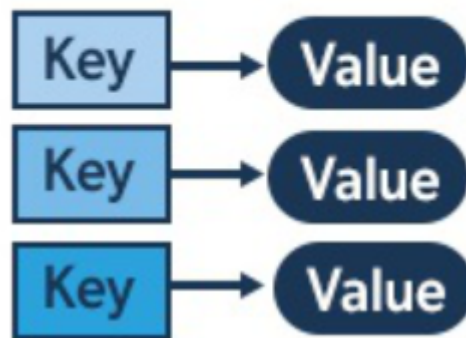


- Colonne
- Document
- Graphe

Modèle clé/valeur

Stockage sous forme d'un dictionnaire clé/valeur

Key-Value



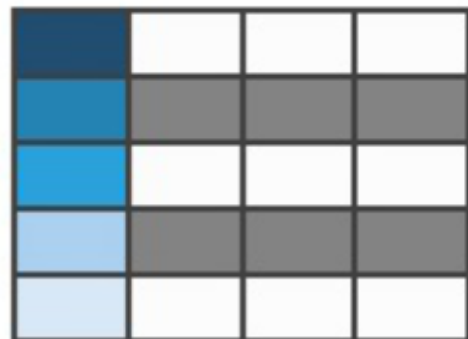
- Clé : chaîne de caractères unique
- Valeur : typage au besoin

ex : redis (StackOverflow), riak (GitHub), Memcached (wikipedia), voldemort (LinkedIn).

Modèle Colonne

Stockage sous forme de table dénormalisée

Column-Family



- système proche des bases de données relationnelles.
- tableau de clé/valeur groupable en famille (+/- table).

ex: Cassandra (Nasa), HBase (Facebook, xiaomi), BigTable (GCP).

Modèle Document

Stockage sous forme de document

Document



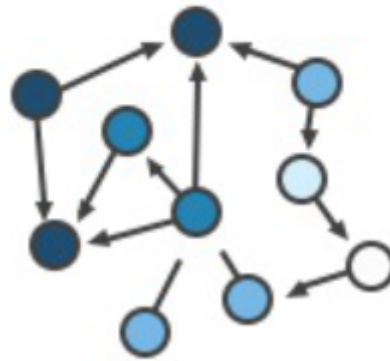
- Object JSON avec un id et des propriétés clé/valeur.
- valeur peuvent etre d'autres documents.

ex: MongoDB (SEGA, ThermoFisher Scientific), CouchDB (CERN)

Modèle Graphe

Stockage sous forme de graphique relationnelle

Graph



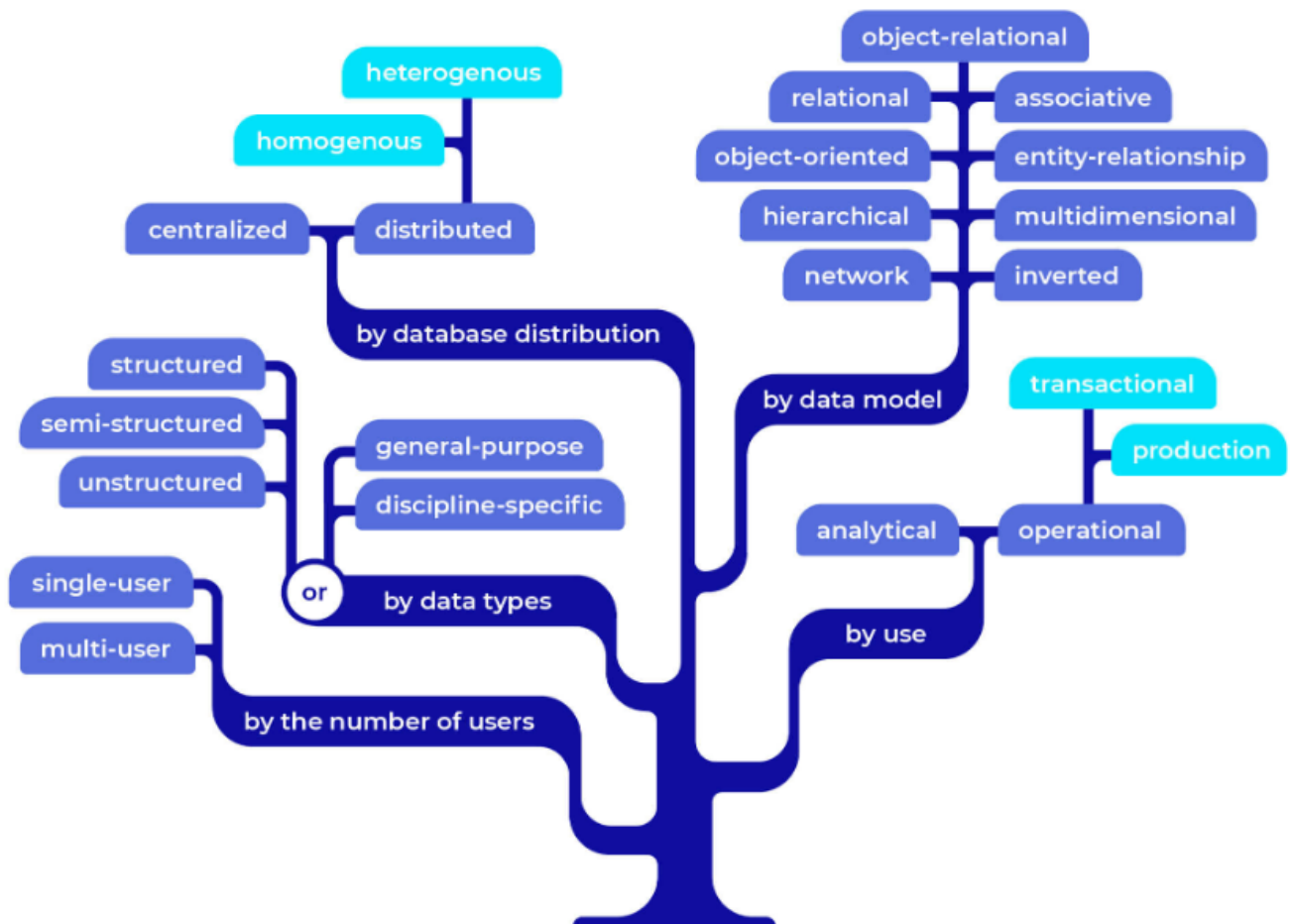
- Noeuds de données de type document clé/valeur
- ARC : relation orientée entre les noeuds porteur de propriétés.

ex: Neo4j (Orange, Airbus), OrientDB, Titan

Choix du modèle de stockage



Critere variable selon besoin



Exercice

Exercice 01

Architectures et algorithmes