

BIG DATA



Les bases de données NoSQL

Introduction

À l'heure du Big Data, les bases de données relationnelles ne sont plus adaptées. Pour prendre en charge les immenses volumes de données, les stocker et les analyser, il est impératif de s'en remettre à de nouvelles solutions.

Une base de données NoSQL est une base de données “non relationnelle”. Il est possible d'y stocker des données sous une forme non structurée, sans suivre de schéma fixe. Les jointures ne sont plus nécessaires, et le scaling est facilité.

Définition

Le NoSQL regroupe de nombreuses bases de données, récentes pour la plupart, qui se caractérisent par une logique de représentation de données non relationnelle et qui n'offrent donc pas une interface de requêtes en SQL.

NoSQL signifie Not OnlySQL et non pas No SQL, il s'agit de compléments aux SGBDR pour des besoins spécifiques et non de solutions de remplacement.

Les bases NoSQL se fondent sur une approche dite *schema-less*, c'est à dire sans schéma logique défini a priori.

Définition

NoSQL est différent de SQL, il ne nécessite pas de schéma et n'a pas de relations de table, il est donc plus flexible.

Les bases de données NoSQL continuent d'augmenter en nombre d'utilisations, en particulier dans les implémentations de données volumineuses et les applications Web en temps réel.

Sa popularité n'a cessé d'augmenter au début de ce siècle millénaire, déclenchée par les besoins des entreprises basées sur le Web 2.0 et des applications gérées.

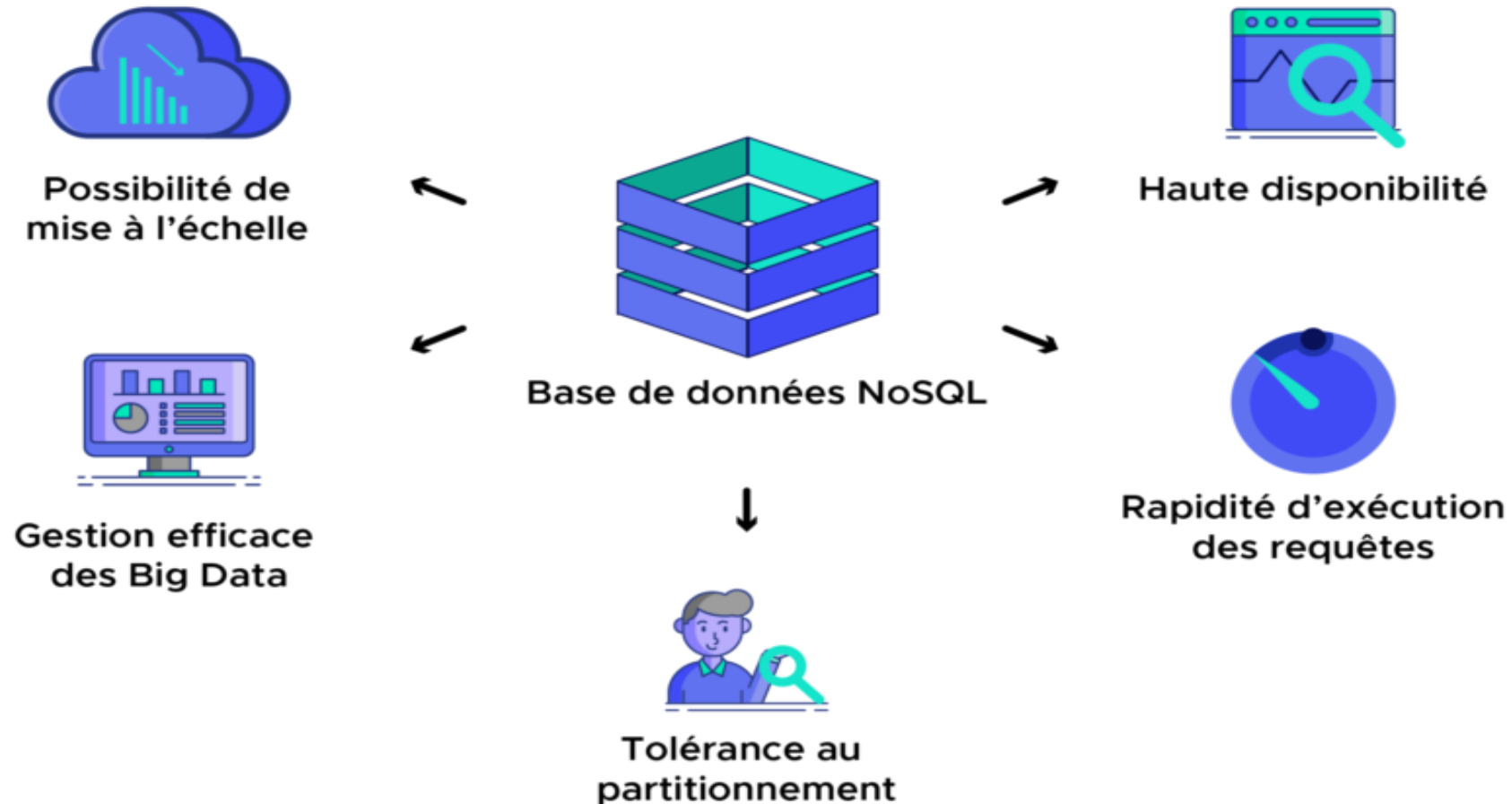
Historique

Le terme NoSQL a été publié pour la première fois par Carlo Strozzi en 1998 pour nommer la base de données qu'il développait « Strozzi NoSQL open-source relational database ». « NoREL », qui fait référence au terme « No Relational ».

Fin 2000, le développement de NoSQL a repris, dans le but de dépasser les limites de SQL, notamment en termes de scalabilité et de potentiel de collecte de données multi-structurées.

Début 2009, Johan Oskarsson, l'un des développeurs de Last.fm, a réintroduit le terme NoSQL lorsqu'il a organisé un événement pour discuter des « bases de données distribuées non relationnelles open source ».

Pourquoi le NoSQL?



De point de vue métier, utiliser un environnement Big Data et NOSQL fournit un avantage compétitif certain

Acteurs du NoSQL

Amazon : DynamoDB, SimpleDB

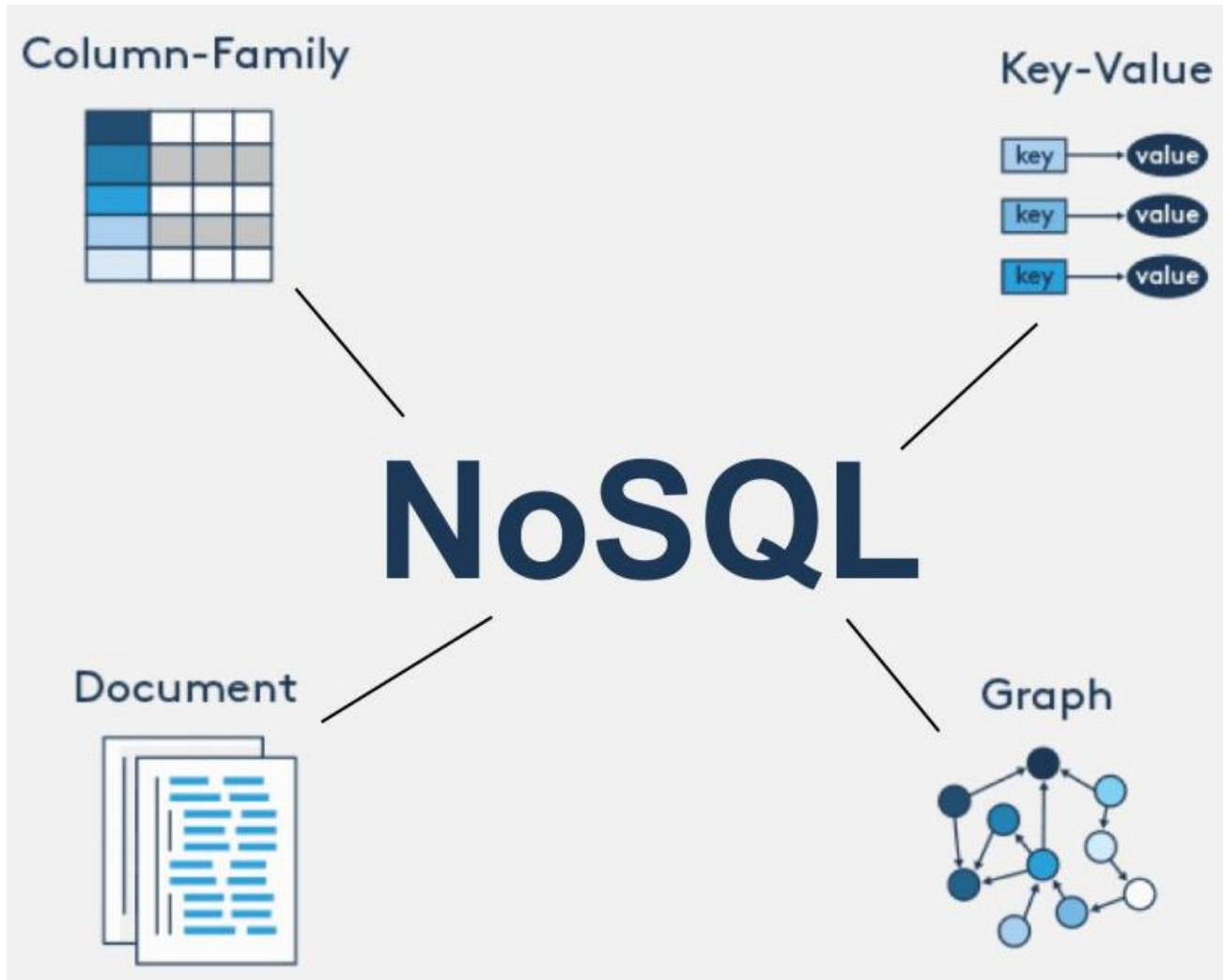
Microsoft : Azure Table Storage

Google : BigTable, Datastore, GFS

Apache : CouchDB, Cassandra, Hadoop/HBase

Beaucoup de start-ups...

Taxonomie des BD NoSQL



Toutefois aucun de ces quatre types de bases de données ne permet de résoudre n'importe quel problème.

Il est nécessaire de **choisir la base de données adéquate** en fonction du cas d'usage.

1/ BD NoSQL Clé-Valeurs

Type basique

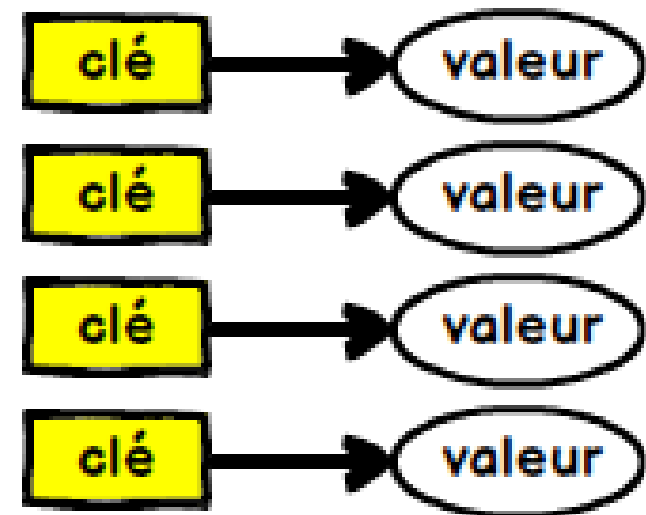
Principes: Représentation des données sous forme de clé/valeur. Les valeurs peuvent être de simple chaînes de caractères ou des objets complexes.

• **Utilisation:** dépôt de données avec besoins de requêtage simples (préférences d'utilisateur, ...)

Exploitation basée sur 4 opérations : Accès par la clé

- **Create:** Création d'un objet.
- **Read:** Lecture d'un objet.
- **Update:** Mise à jour d'un objet.
- **Delete:** Suppression un objet.

Clé-Valeur



2/ Base de données orientées document

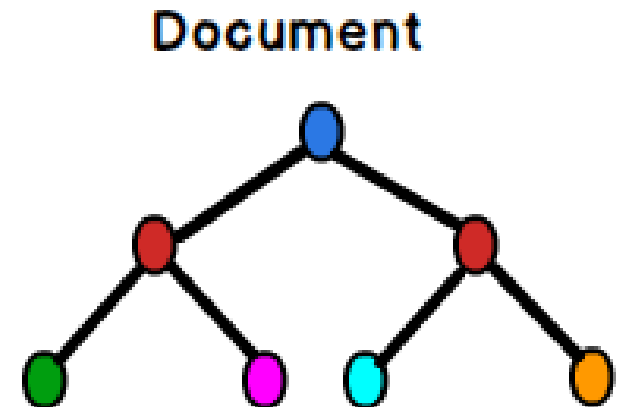
Principes: c'est une variante des SGBD clé/valeur, où la valeur est un document de type XML ou JSON.

–Les documents ont une structure arborescente, ils sont composés de champs et des valeurs associées.

Ce type de SGBD permet d'effectuer des requêtes sur le contenu des documents. Ce type de base de données offre une flexibilité accrue.

Utilisé pour les systèmes CMS, les plateformes de blogging, ou les applications de e-commerce

Ne convient pas pour les transactions complexes nécessitant des opérations ou des requêtes multiples sur des structures agrégées variables.



3/ Base de données orientées colonnes

Principes: Repose sur des colonnes. Proche du relationnel. Mais le stockage des données se fait par colonne et non par ligne.

Chaque colonne est traitée séparément, et les valeurs sont stockées de façon contigüe.

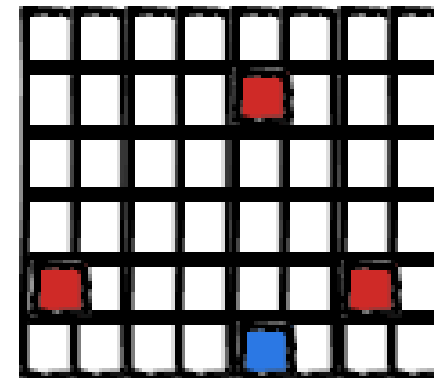
Ajout de colonnes facile et dynamique.

Possibilité de compression des données

• **Utilisation:** Cette catégorie de base de données offre de hautes performances pour les requêtes d'agrégation comme SUM, COUNT, AVG et MIN. Pour cause les données sont déjà disponibles et prêtes dans une colonne.

Orienté colonne

APACHE
HBASE

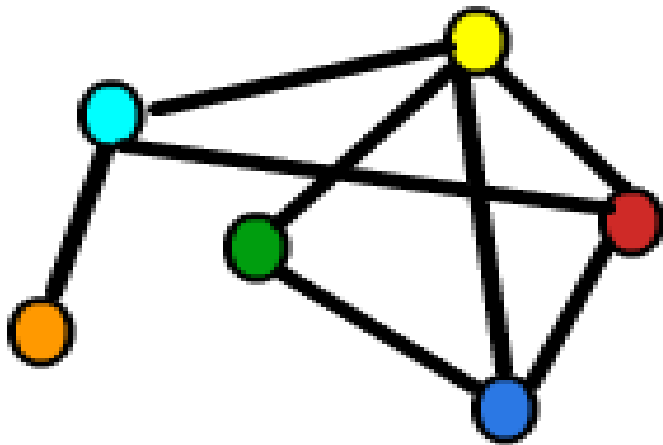


4/ Base de données orientées graphe

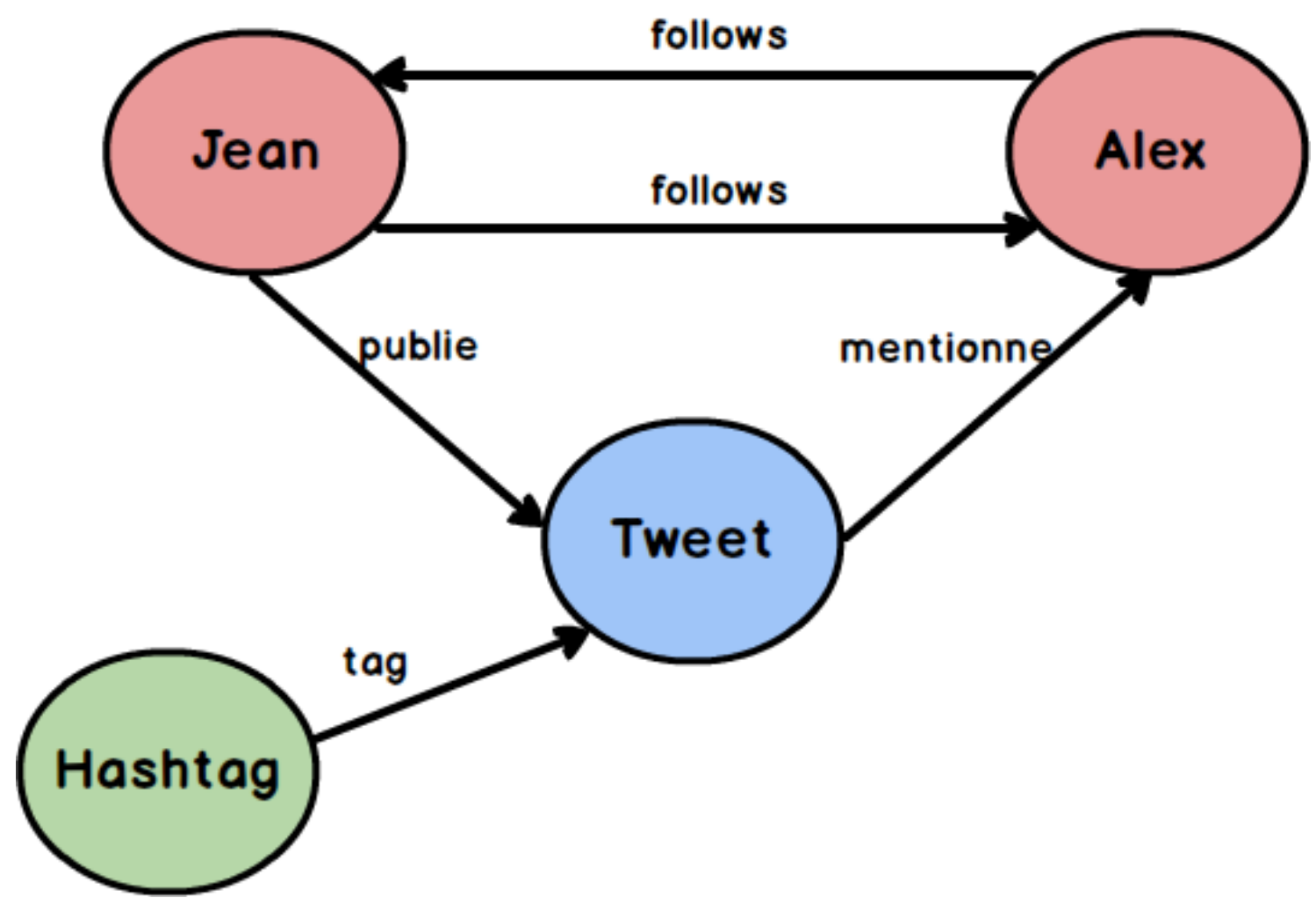
Principes:

- Les données sont représentées sous-forme de graphe : Des nœuds pour les entités, des arcs pour les relations entre les entités.
- Ce type de SGBD est adapté à la manipulation de données fortement connectées
- **Utilisation:** Systèmes de recommandations, Réseaux sociaux, Systèmes de transport . .

Graphe



4/ Base de données orientées graphe



Exemple : HBase

HBase est un système de stockage efficace pour des données très volumineuses
Modèle de données orienté colonne .

Il permet d'accéder aux données très rapidement même quand elles sont gigantesques.

HBase est utilisée par FaceBook pour stocker tous les messages SMS, email et chat...

HBase peut être utilisée à la fois comme:

–Base de données temps réel.

–Base de données pour une lecture intensive pour les systèmes décisionnels.

Clés

isbn7615

isbn7615

isbn7892

isbn7892

Colonnes et Valeurs

colonne=auteur valeur="Jules Verne"

colonne=titre valeur="De la Terre à la Lune"

colonne=auteur valeur="Jules Verne"

colonne=titre valeur="Autour de la Lune"

Exemple : HBase

Pour obtenir une grande efficacité, les données des tables Hbase sont séparées en **régions**.

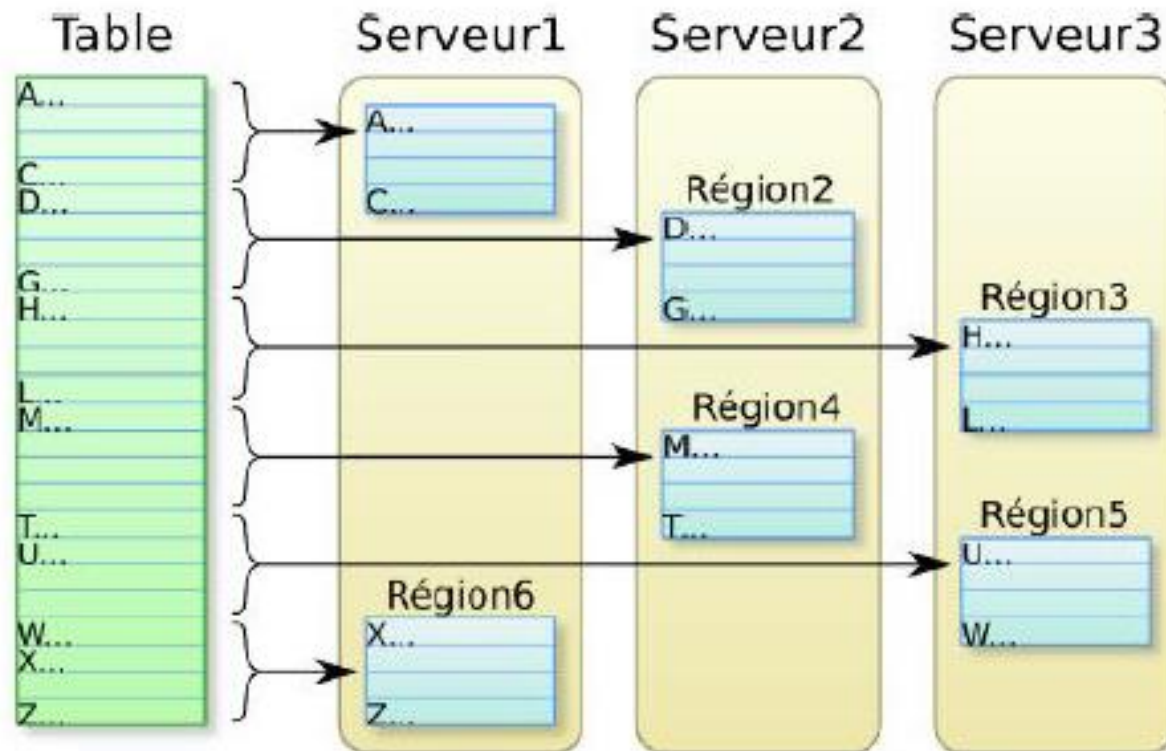
- Une région contient un certain nombre de n-uplets contigus (un intervalle de clés successives).
- Lorsqu'on crée une table, elle est mise dans **une seule région**.
- Lorsque la table dépasse une certaine limite, elle se fait couper en deux régions au milieu de ses clés. Et ainsi de suite si les régions deviennent trop grosses.
- Chaque région est gérée par un **Serveur de Région** (*RegionServer*).
- Ces serveurs sont distribués sur le cluster, ex: un par machine.
- Un même serveur de région peut s'occuper de plusieurs régions de la même table.

Au final, les données sont stockées sur HDFS.

Exemple : HBase

Une table est découpée en régions faisant à peu près la même taille.

- Le découpage est basé sur les clés.
- Chaque région est gérée par un Serveur de région.
- Un même serveur peut gérer plusieurs régions.



Comparaison

Modèle	Performance	Évolutivité	Flexibilité	Complexité
Clef/Valeur				
Orienté Colonnes				
Orienté Document				
Orienté Graphes				

Illustrations

Représentation de ventes en relationnel

Table Sales

#ticket	#date	#book
1	01/01/16	2212121504
1	01/01/16	2212141556
2	01/01/16	2212141556

Table Book

#isbn	#title	#author
2212121504	Scenari	1
2212141556	NoSQL	2

Table Author

#id	surname	firstname
1		
2	Bruchez	Rudi



Représentation de ventes en colonne

Family Sales

#ticket	date	books
1	01/01/16	2212121504 2212141556
2	01/01/16	2212141556

Family Book

#isbn	title	a-surname	a-firstname
2212121504	Scenari		
2212141556	NoSQL	Bruchez	Rudi

Illustrations

Représentation de ventes en relationnel

Table Sales

#ticket	#date	#book
1	01/01/16	2212121504
1	01/01/16	2212141556
2	01/01/16	2212141556

Table Book

#isbn	#title	#author
2212121504	Scenari	1
2212141556	NoSQL	2

Table Author

#id	surname	firstname
1		
2	Bruchez	Rudi



Représentation de ventes en document

Collection Sales

#oid	
4d040766076 6b236450b45 a3	"ticket" : 1 "date" : "01/01/16" "books" : [<ul style="list-style-type: none"> "_id" : 2212121504 "_id" : 2212141556]
4d040766076 6b236450b45 a4	"ticket" : 2 "date" : "01/01/16" "books" : [<ul style="list-style-type: none"> "_id" : 2212141556]

Collection Book

#oid	
4d040766076 6b236450b45 a5	"isbn" : 2212121504 "title" : "Scenari"
4d040766076 6b236450b45 a6	"isbn" : 2212141556 "title" : "NoSQL" "author" : { <ul style="list-style-type: none"> "surname" : Bruchez "firstname" : Rudi }

Illustrations

Représentation de ventes en relationnel

Table Sales

#ticket	#date	#book
1	01/01/16	2212121504
1	01/01/16	2212141556
2	01/01/16	2212141556

Table Book

#isbn	#title	#author
2212121504	Scenari	1
2212141556	NoSQL	2

Table Author

#id	surname	firstname
1		
2	Bruchez	Rudi



Représentation de ventes en graphe

Classe Sales

#oid		
4d040766076 6b236450b45 a3	<i>property</i>	ticket : 1
	<i>property</i>	date : 01/01/16
	<i>relation</i>	book : 4d0407660766b236450b45a5
	<i>relation</i>	book : 4d0407660766b236450b45a6
4d040766076 6b236450b45 a4	<i>property</i>	ticket : 2
	<i>property</i>	date : 01/01/16
	<i>relation</i>	book : 4d0407660766b236450b45a6

Classe Book

#oid		
4d040766076 6b236450b45 a5	<i>property</i>	title : Scenari
4d040766076 6b236450b45 a6	<i>property</i>	title : NoSQL
	<i>relation</i>	author : 4d0407660766b236450b45a8

Classe Author

#oid		
4d040766076 6b236450b45 a8	<i>property</i>	sumame : Bruchez
	<i>property</i>	firstnam : Rudi

BDD SQL VS BDD NoSql

Différence	SQL	NoSQL
Définition	SGBDR ou bases de données relationnelles	Base de données non relationnelle ou base de <i>données distribuée</i>
Utilisation	Requête pour analyser et récupérer données	Traiter des données liées à des applications et des sites Web modernes de plus en plus complexes
Langage de requête	Langage de requête structuré (SQL)	Ne nécessite pas un langage de requête trop complexe
Structure de la base de données	En forme de table	Valeur-clé, d'une colonne, d'un document et d'un graphique
Schéma	Besoin d'être déterminé d'abord	Schéma dynamique pour les données non structurées
Adapté pour	Requêtes complexes et intensives	Grande base de données, Big Data
Exemples	MySQL, Postgres, MS-SQL	Redis, Neo4j, MongoDB