

CentraleSupélec  

Big Data : Informatique pour les données et calculs massifs

9 – Technologie des moteurs de Bdd NoSQL

Stéphane Vialle

université CentraleSupélec PARIS-SACLAY Sciences et technologies de l'information et de la communication (STIC)  

Stephane.Vialle@centralesupelec.fr
<http://www.metz.supelec.fr/~vialle>

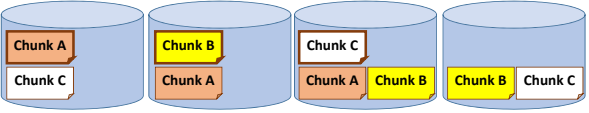
CentraleSupélec

Concepts de sharding et répliquions

2

CentraleSupélec

Principe de *sharding* et répliquion



Distribution et répliquion:

- Un fichier/une table est découpé en morceaux (*chunks*) distribués pour permettre:
 - des accès parallèles plus rapides (répartir la charge)
 - des stockages plus volumineux (et extensibles)
- Les *chunks* sont répliqués sur des nœuds différents pour la tolérance aux pannes
- Le *sharding* est un découpage des différentes lignes d'une table par tranches dans différents *chunks* (au lieu de distribuer les colonnes)

Concepts d'HDPS

Principe de Sharding et réplication

Des services de :

- Cartographie du stockage
- Monitoring & Hearbeat

Réplication des services pour la tolérance aux pannes

Concepts d'HDFS

Architecture de HBASE
(Microsoft & OpenSrc)

5

Archi. HBase au dessus d'Hadoop

BdD orientée colonnes

100	vente-2010	100000	vente-2011	150000	vente-2014	180000
200	vente-2012	1000				
211	vente-2010	500000	achat-2016	10000		

- Une ligne est indexée par une clé unique
- Les colonnes sont regroupées en familles
Une famille est stockée dans un même fichier *Hfile*
- Des options de compression peuvent être associées à une famille de colonnes
- Un *versioning* est disponible pour chaque colonne
- Bâtie au dessus d'Hadoop et d'HDFS

Archi. HBase au dessus d'Hadoop

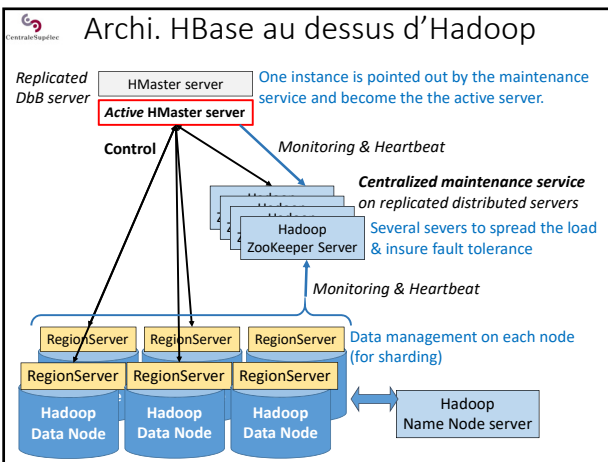
BdD orientée colonnes

2 niveaux de gestion de la distribution des données :

- *Sharding* des tables de données en *Hfiles*
→ fait par des *Serveurs de Région* de Hbase (1 par nœud)
- Gestion/distribution/réplication des fichiers *Hfiles*
→ fait par *HDFS d'Hadoop*

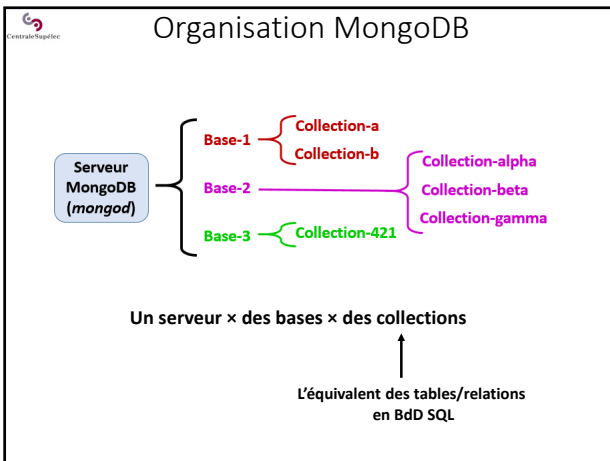
Sharding progressif d'une table :

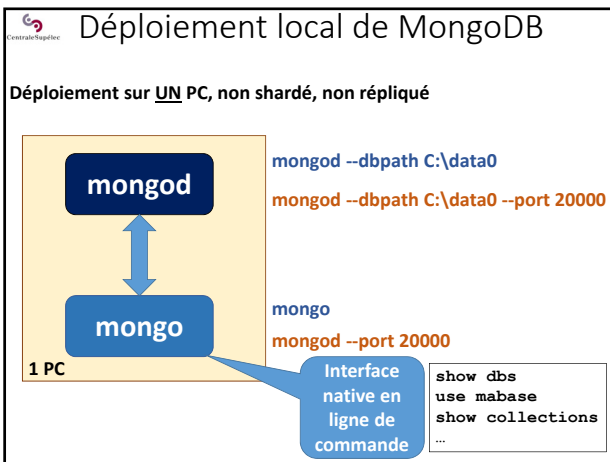
1. Au début une table est contenue dans une région (un nœud)
2. Quand elle grossit, une région est scindée en 2 sous-régions sur le même serveur de région (le même nœud)
3. Séparation signalée au maître (le *Hmaster*)
4. Peut demander la migration d'une sous-région sur une autre serveur de région

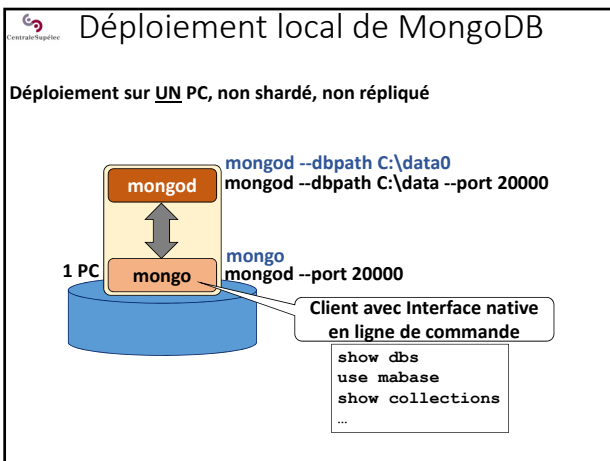


Architecture et *mapReduce* de MongoDB

9







mapReduce de MongoDB (2)

Principes du **mapReduce** de MongoDB :

input data → query() → map() → Shuffle & Sort → reduce() → finalize() → output data

- **map()** fonctionne comme en Hadoop, mais **sur une seule collection**
- pour éviter d'implanter un filtrage en JS, une **query** exprimée en Mongo shell est applicable en amont (plus pratique et plus rapide)
- **reduce()** est **plus contraint qu'en Hadoop** car MongoDB applique **reduce()** sur des **< key , [val] >** de sortie de map(), et sur des sorties précédentes de **reduce()** :
 - **le format de sortie de reduce() doit être celui de sortie de map()**
 - **La fonction reduce() doit être commutative, associative et idempotente**
- **finalize()** permet de modifier la sortie finale de **reduce()** sans contrainte

16

mapReduce de MongoDB (3)

Principes du **mapReduce** de MongoDB :

input data → query() → map() → Shuffle & Sort → reduce() → finalize() → output data

- **Pas de Combiner**
- **Pas d'optimisation** possible du **Shuffle & Sort**

Les contraintes sur la fonction **reduce()** :

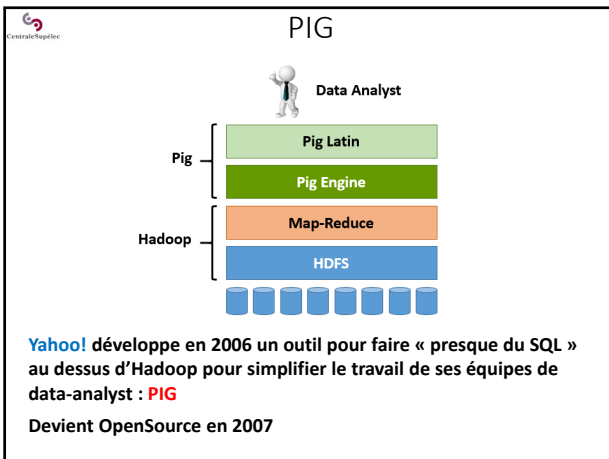
- **évitent les débordements mémoire**, en traitant chaque paire < clé, liste de val > en plusieurs passes si la liste de valeurs est « grande »
- **mais...limitent fortement les algorithmes possibles** en une seule opération **mapReduce**.

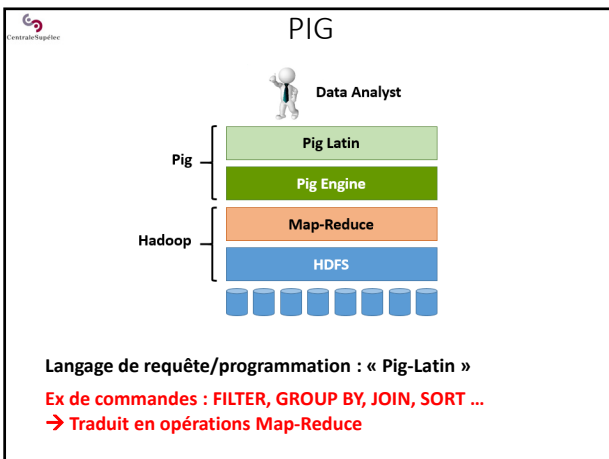
17

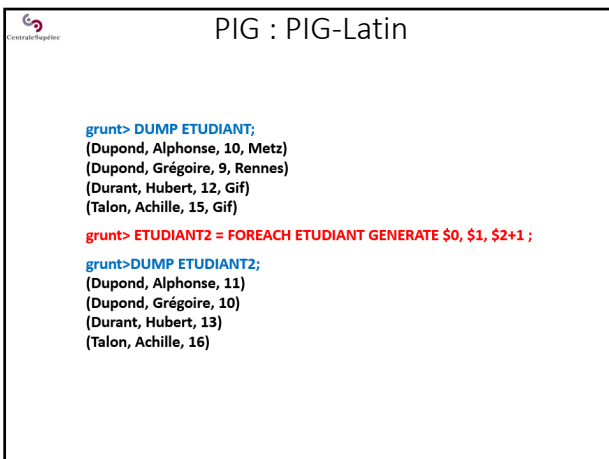
Environnements NoSQL de plus haut niveau


- PIG
- HIVE

18







 PIG : PIG-Latin

```


grunt> DUMP VEHICULE;
(AA 123 AB, bleue, Twingo)
(AZ 451 GT, noire, C3)
(BC 634 FY, jaune, Twingo)
(DE 398 AA, blanche, DS4)

grunt> DUMP CONSTRUCTEUR;
(Twingo, Renault)
(C3, Citroen)
(DS4, Citroen)

grunt> VOITURE = JOIN VEHICULE BY $2, CONSTRUCTEUR BY $0;

grunt> DUMP VOITURE;
(AA 123 AB, bleue, Twingo, Renault)
(AZ 451 GT, noire, C3, Citroen)
(BC 634 FY, jaune, Twingo, Twingo, Renault)
(DE 398 AA, blanche, DS4, DS4, Citroen)

```


 PIG : PIG-Latin

```

grunt> VOITURE2 = FOREACH VOITURE GENERATE $0, $1, $2, $4 ;

grunt> VOITURE3 = ORDER VOITURE2 BY $2, $0 DESC ;
(AZ 451 GT, noire, C3, Citroen)
(DE 398 AA, blanche, DS4, Citroen)
(BC 634 FY, jaune, Twingo, Renault)
(AA 123 AB, bleue, Twingo, Renault)

```


 PIG : PIG-Latin

Stratégie :

- Traduit les requêtes en **DAG d'opérations Map-Reduce** (!!)
- Exprime un maximum de parallélisme entre les opérations du DAG (souci de performance)
- Le « PIG-engine » exécute le DAG, et fixe le nombre de *Reducers au mieux* (mais on peut le guider)

Rmq :

- **SQL : langage déclaratif** : on dit ce que l'on veut obtenir
- **Pig-Latin : langage procédural** : on dit quoi faire (pour obtenir le résultat).

 Hive

Facebook 2005 : encore plus proche de SQL que PIG

- Hive : langage déclaratif proche de SQL
- Pour des data-analysts qui ne peuvent pas programmer du Map-Reduce (pas programmer en langage procédural)
- Génère du Map-Reduce au dessus d'Hadoop.

 Technologie des moteurs de Bdd NoSQL