

CentraleSupélec

Big Data : Informatique pour les données et calculs massifs

## 9 – Technologie des moteurs de Bdd NoSQL

Stéphane Vialle

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

CentraleSupélec

## Concepts de sharding et répliquions

2

CentraleSupélec

### Principe de *sharding* et répliation

**Distribution et répliation:**

- 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'HDFS

CentraleSupélec

### Principe de Sharding et répliation

**Des services de :**

- Cartographie du stockage
- Monitoring & Hearbeat

**Répliation des services pour la tolérance aux pannes**

Concepts d'HDFS

CentraleSupélec

## Architecture de HBASE (Microsoft & OpenSrc)

5

CentraleSupélec

## 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 *versionning* est disponible pour chaque colonne
- Bâtie au dessus d'Hadoop et d'HDFS

### Archi. HBase au dessus d'Hadoop

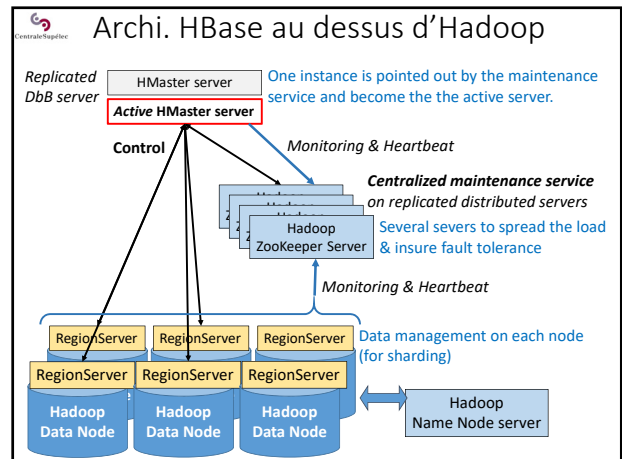
**BdD orientée colonnes**

**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*

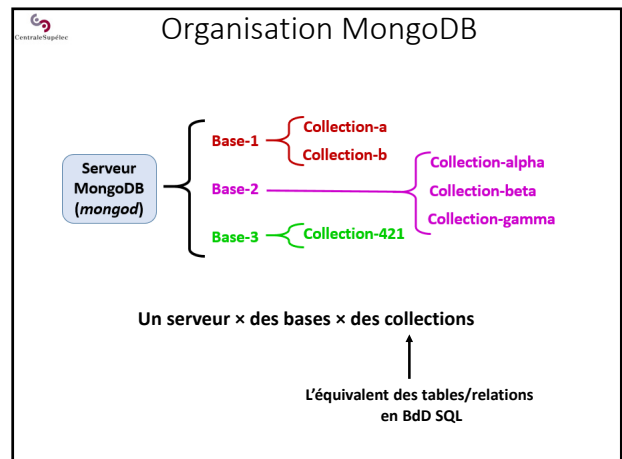
**Sharding progressif d'une table :**

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



### Déploiement local de MongoDB

Déploiement sur UN PC, non shardé, non répliqué

mongod --dbpath C:\data0  
mongo --dbpath C:\data0 --port 20000  
mongo  
mongo --port 20000

Interface native en ligne de commande

```
show dbs
use mabase
show collections
...
```

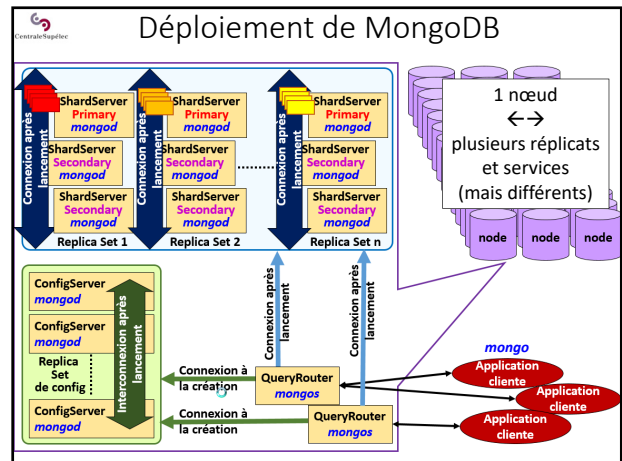
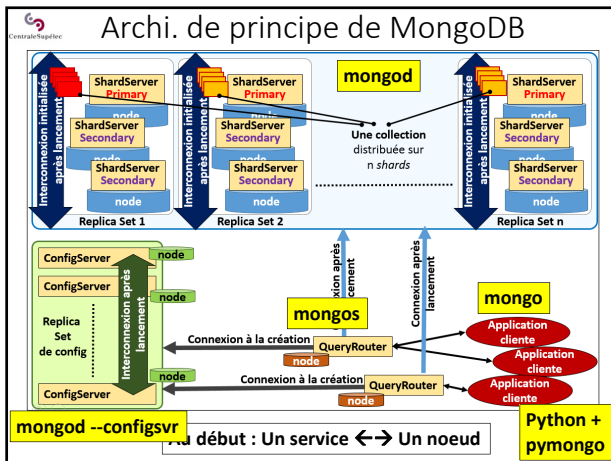
### Déploiement local de MongoDB

Déploiement sur UN PC, non shardé, non répliqué

mongod --dbpath C:\data0  
mongo --dbpath C:\data --port 20000  
mongo --port 20000

Client avec Interface native en ligne de commande

```
show dbs
use mabase
show collections
...
```



### mapReduce de MongoDB (1)

**MongoDB possède son propre « Map-Reduce » :**

- Ses propres règles de fonctionnement (un peu différentes d'Hadoop)
- Ses propres règles de déploiement des tâches
- Et son propre middleware sous-jacent (il n'est pas bâti au dessus d'Hadoop)
- Fonctionne sur des bases distribuées (*sharded*)

**Principes du mapReduce de MongoDB :**

- Une *query* pour pré-filtrer la collection traitée
- Une fonction *map()*, en JavaScript et qui accède à la base
- Une fonction *reduce()*, en JavaScript et qui ne doit PAS accéder à la base (seulement aux résultats du *map()*), qui doit être *commutative, associative et idempotente* (!!)
- Une fonction *finalize()*, en JavaScript et qui ne doit pas accéder à la base
- La possibilité de définir un ensemble de variables globales aux 3 fonctions *map()*, *reduce()* et *finalize()*

### mapReduce de MongoDB (2)

**Principes du mapReduce de MongoDB :**

The flowchart shows the MapReduce process: 'query()' leads to 'map()', which leads to 'Shuffle & Sort', then to 'reduce()', 'finalize()', and finally 'output data'. A note indicates 'Identical < key, val > scheme'.

- *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

### mapReduce de MongoDB (3)

**Principes du mapReduce de MongoDB :**

The flowchart shows the MapReduce process with a lock icon on the 'Shuffle & Sort' step. The flow is: 'query()' -> 'map()' -> 'Shuffle & Sort' -> 'reduce()' -> 'finalize()' -> 'output data'. A note indicates 'Identical < key, val > scheme'.

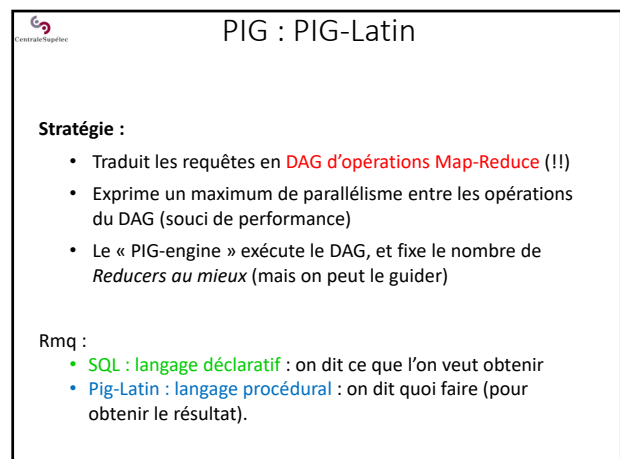
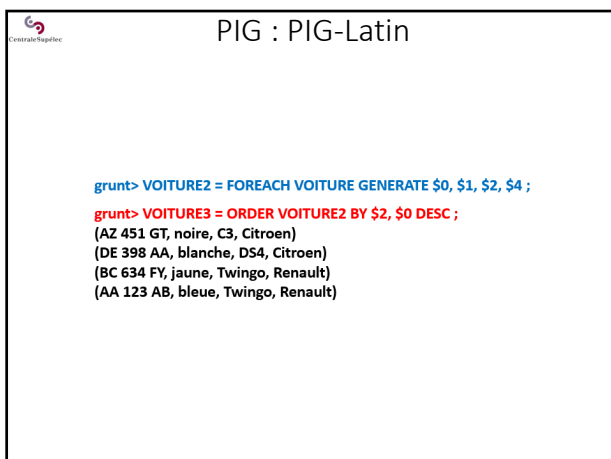
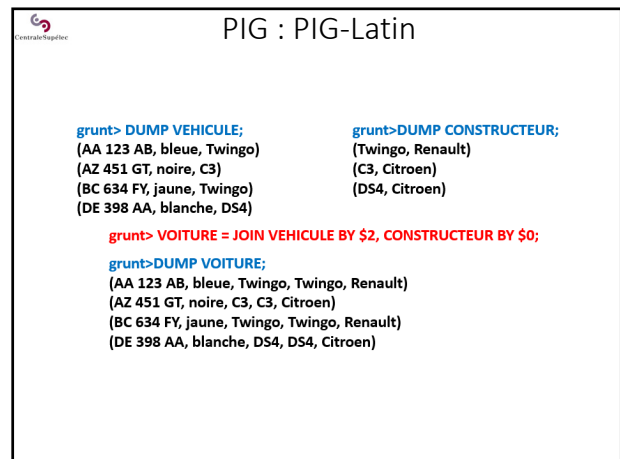
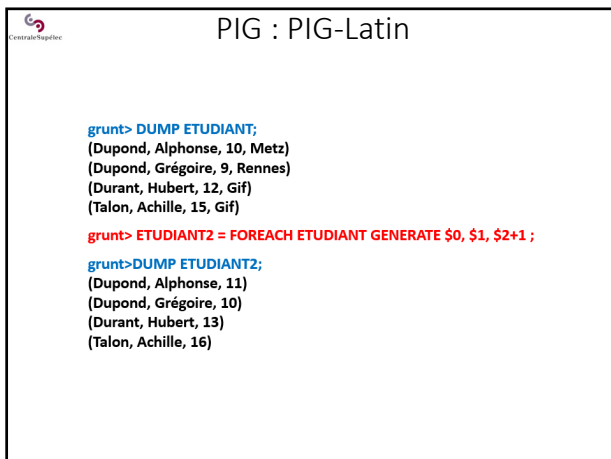
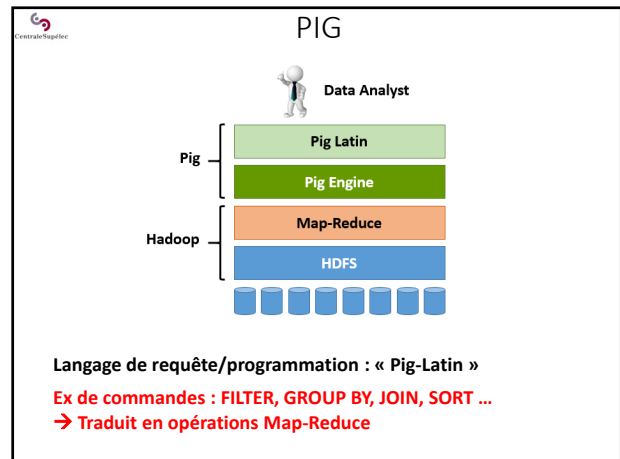
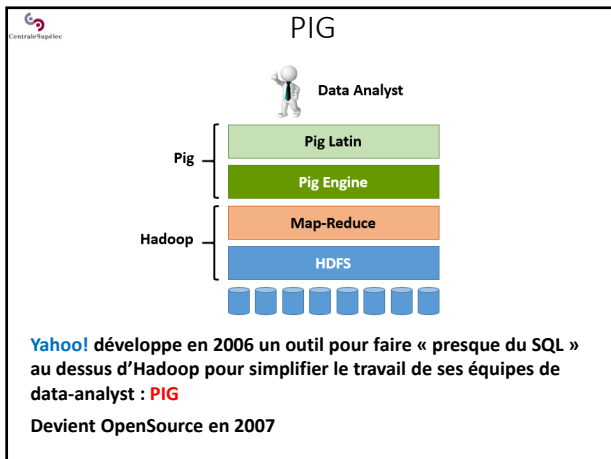
- Pas de *Combiner*
- Pas d'optimisation possible du *Shuffle & Sort*


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

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

### Environnements NoSQL de plus haut niveau

- PIG
- HIVE



 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

