




CentraleSupélec

Big Data : Informatique pour les données et calculs massifs


6 – Algorithmique Map-Reduce

Stéphane Vialle




universit  PARIS-SACLAY

 COLE DOCTORALE
Sciences et technologies
de l'information
et de la communication (STIC)




RISEGrid
Research Institute for Sustainable Electric Grids



Grand Est
ALSACE CHAMPAGNE-ARDENNE LORRAINE

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



CentraleSupélec

D marche algorithmique

- **D finir des t ches *Mappers* et *Reducers***
 - B tir la base de la solution *Map-Reduce*
- **Identifier des routines *Combiner***
 - Eviter des E/S temporaires sur disque (*Hadoop*)
 - R duire le trafic sur le r seau entre *Mappers* et *Reducers*
 - Eviter de saturer la m moire des *Reducers*
- **Identifier le nombre de *Reducers* optimal**
 - R partir la charge au maximum entre plusieurs machines
- **Concevoir des *Partitioners* pour optimiser la solution**
 - Am liorer la r partition de charge des *Reducers*
(si on connait la vari t  des cl s et la distribution des paires)
- **Ecrire des fonctions *keyComparator()* et *groupComparator()***
 - Pour profiter des fonctions de tri int gr es au *Shuffle & Sort*

➔ « Patterns » ou « Patrons » de conception types



Familles de patrons Map-Reduce

Patrons de récapitulation (*Summarization*)

Patrons de filtrage

Patrons de restructuration des données

Patrons de tris

Patrons de jointure

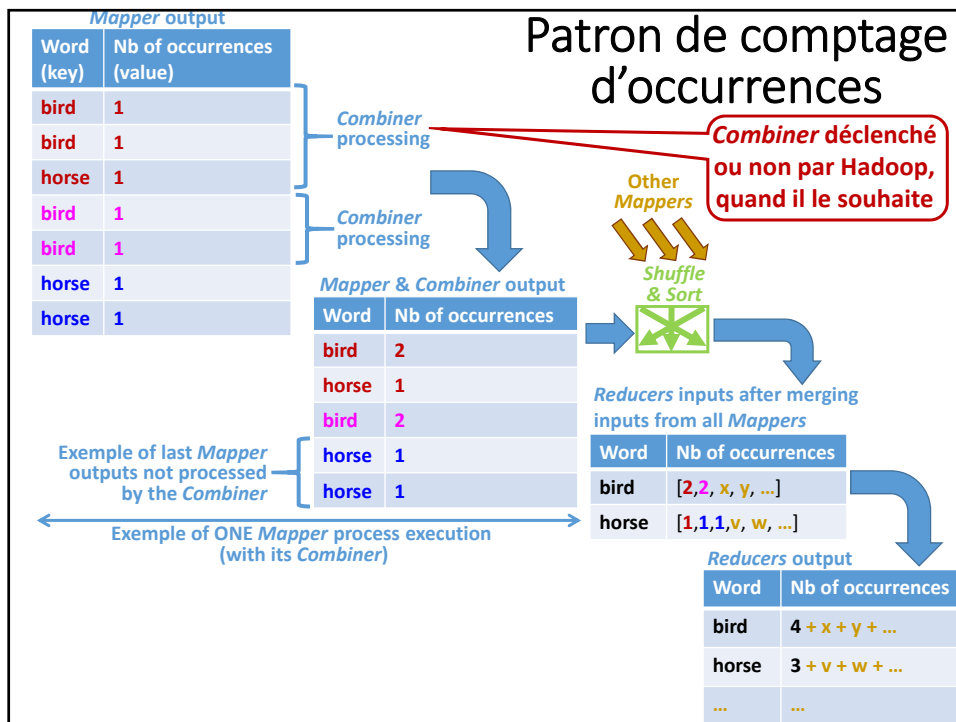
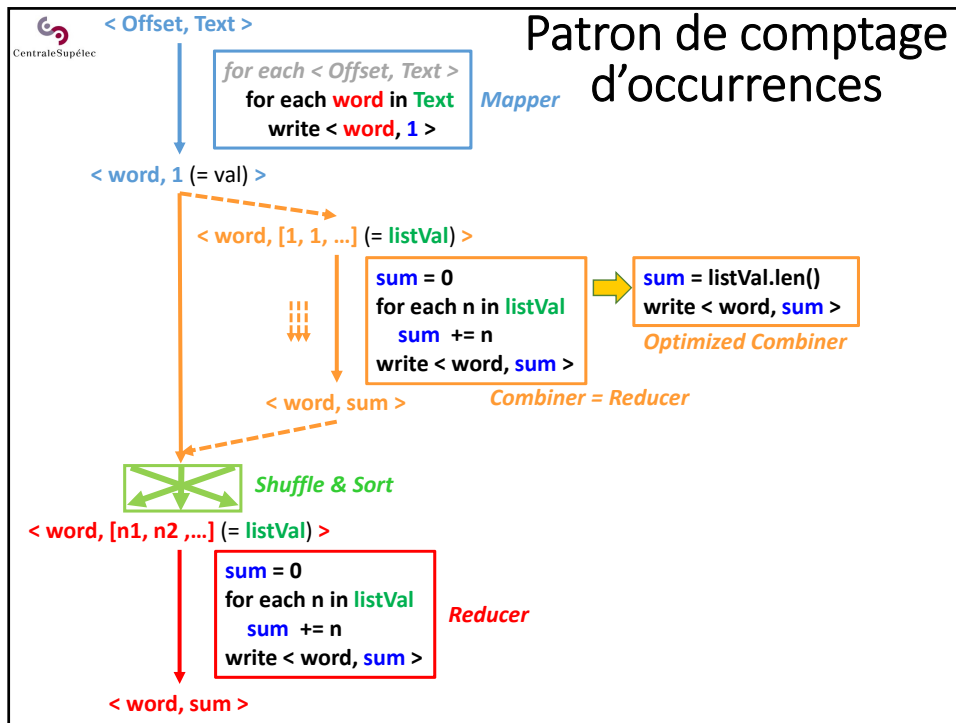


Patrons de récapitulation

1 – Patron de comptage d'occurrences de termes

2 – Patron de calcul d'une moyenne

3 – Patron de réalisation d'un index inversé



CentraleSupélec

Patron de comptage d'occurrences

Déploiement :

```

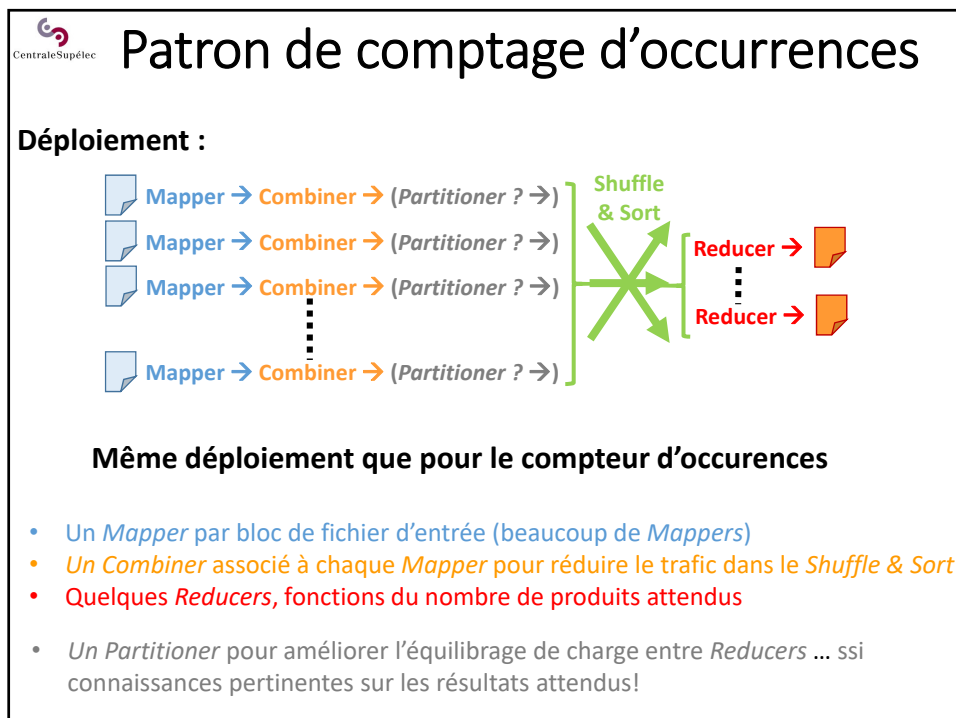
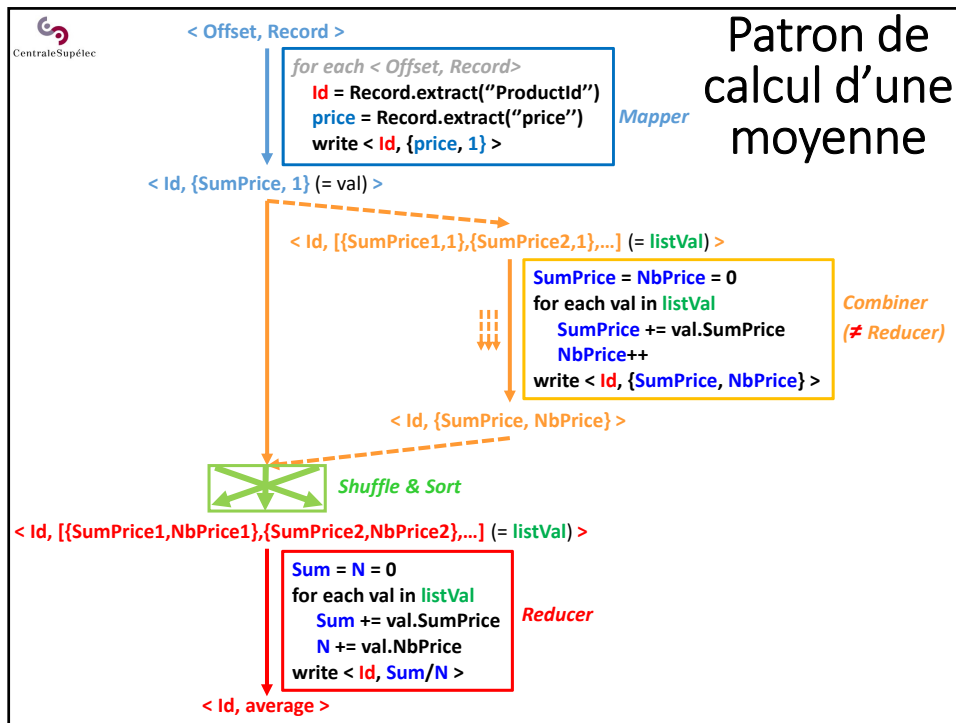
graph LR
    subgraph Mappers
        M1[Mapper] --> C1[Combiner]
        M2[Mapper] --> C2[Combiner]
        M3[Mapper] --> C3[Combiner]
        M4[Mapper] --> C4[Combiner]
    end
    C1 --> PS1["(Partitioner ? →)"]
    C2 --> PS2["(Partitioner ? →)"]
    C3 --> PS3["(Partitioner ? →)"]
    C4 --> PS4["(Partitioner ? →)"]
    PS1 --> SS["Shuffle & Sort"]
    PS2 --> SS
    PS3 --> SS
    PS4 --> SS
    SS --> R1[Reducer]
    SS --> R2[Reducer]
    
```

- Un *Mapper* par bloc de fichier d'entrée (beaucoup de *Mappers*)
- Un *Combiner* associé à chaque *Mapper* pour réduire le trafic dans le *Shuffle & Sort*
- Quelques *Reducers*, fonctions du nombre de termes attendus
- Un *Partitioner* pour améliorer l'équilibrage de charge entre *Reducers* ... ssi connaissances pertinentes sur les résultats attendus!

CentraleSupélec

Patrons de récapitulation

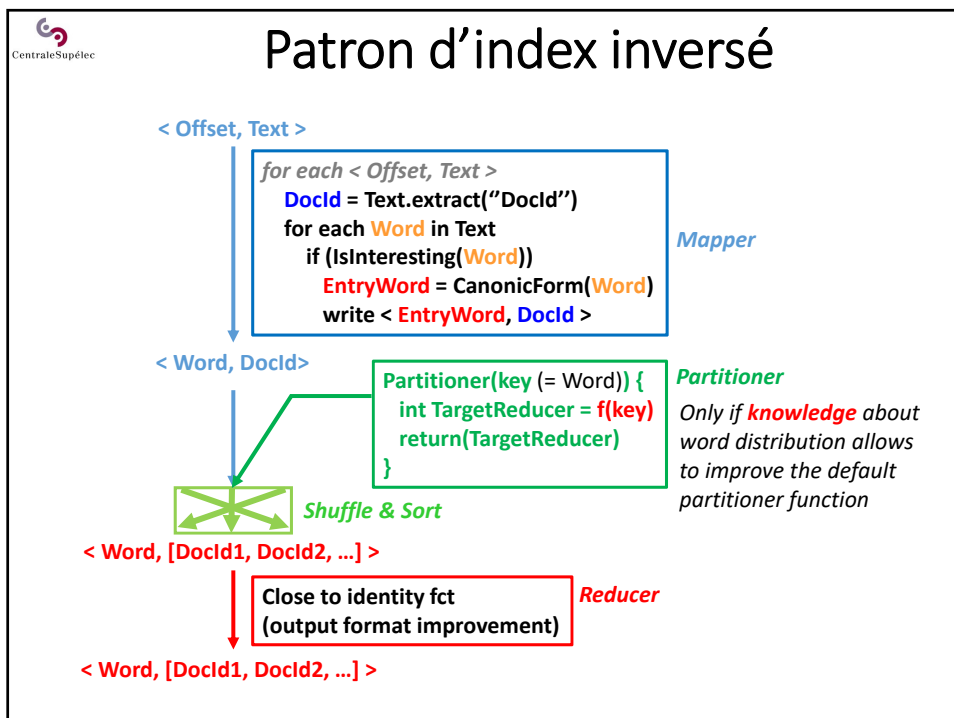
- 1 – Patron de comptage d'occurrences de termes
- 2 – Patron de calcul d'une moyenne
- 3 – Patron de réalisation d'un index inversé



CentraleSupélec

Patrons de récapitulation

- 1 – Patron de comptage d'occurrences de termes
- 2 – Patron de calcul d'une moyenne
- 3 – Patron de réalisation d'un index inversé



CentraleSupélec

Patron d'index inversé

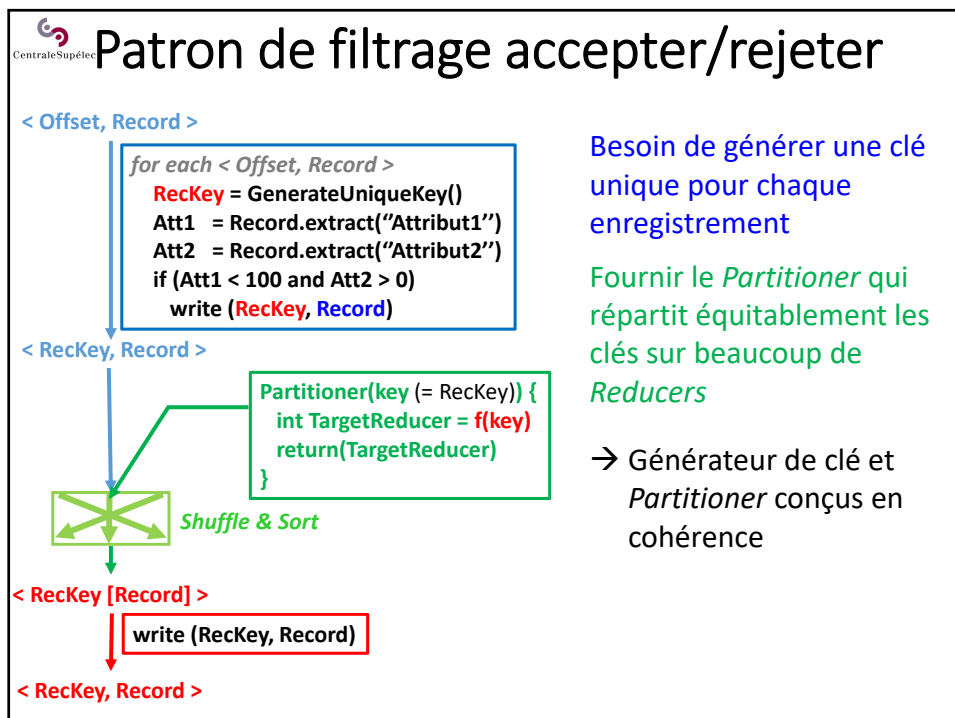
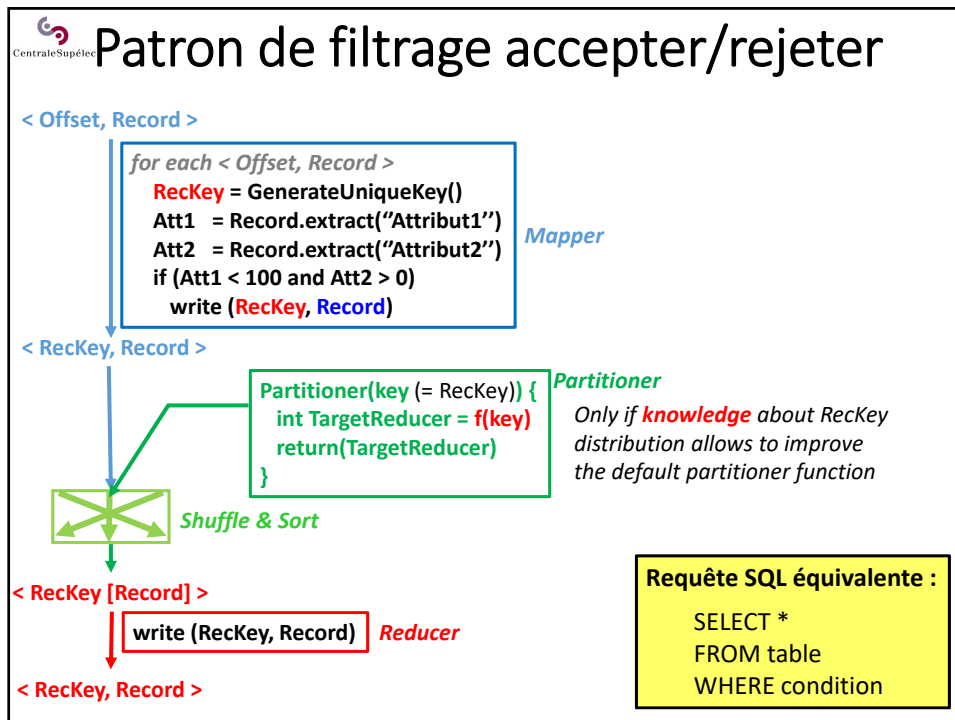
Déploiement :

- *Combiner* inutile : non implémenté
- *Partitioner* pour améliorer l'équilibrage de charge entre *Reducers*
 - informations sur la distribution des mots ?
 - en fonction de la langue des documents ?
 - calcul d'une fonction de hashage optimale

CentraleSupélec

Patrons de Filtrage

- 1 – Patron de filtrage accepter/rejeter
- 2 – Patron de filtrage *Top Ten*
- 3 – Patron de filtrage dé-duplicateur



CentraleSupélec

Patron de filtrage accepter/rejeter

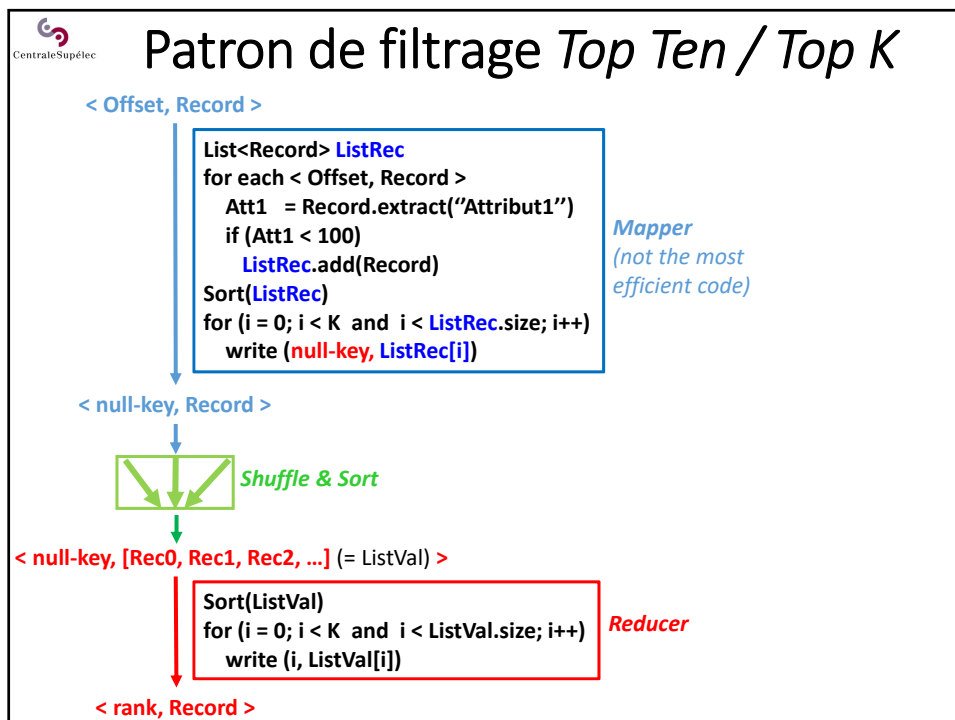
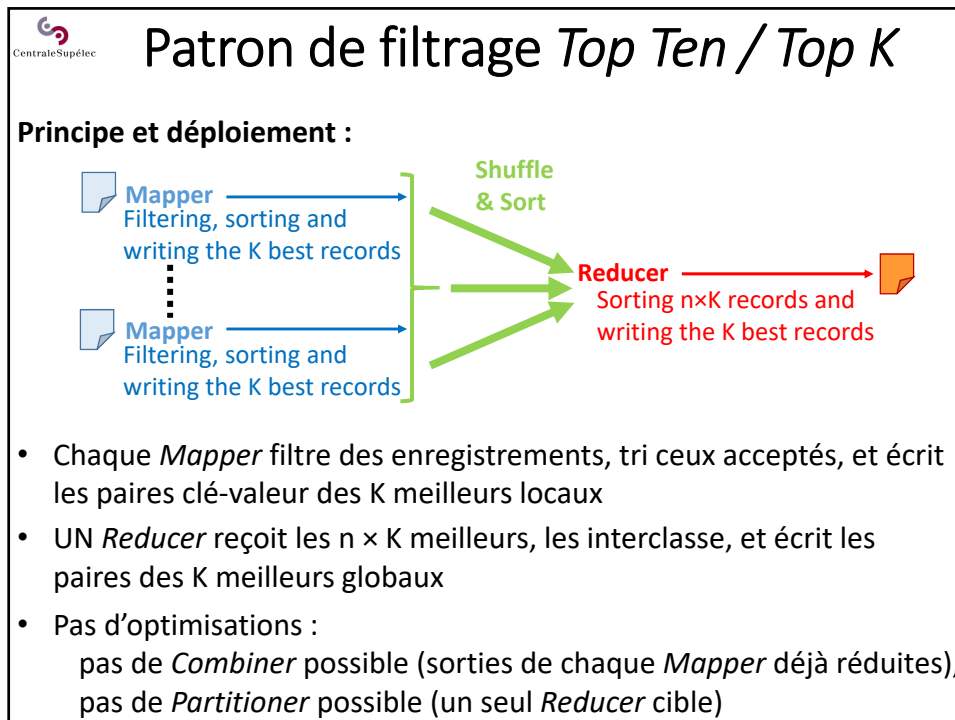
Déploiement :

- *Combiner* inutile (pas de réduction de trafic) : non implémenté
- *Partitioner* pour améliorer l'équilibrage de charge entre *Reducers*
→ clés générées et *Partitioner* conçus ensemble

CentraleSupélec

Patrons de Filtrage

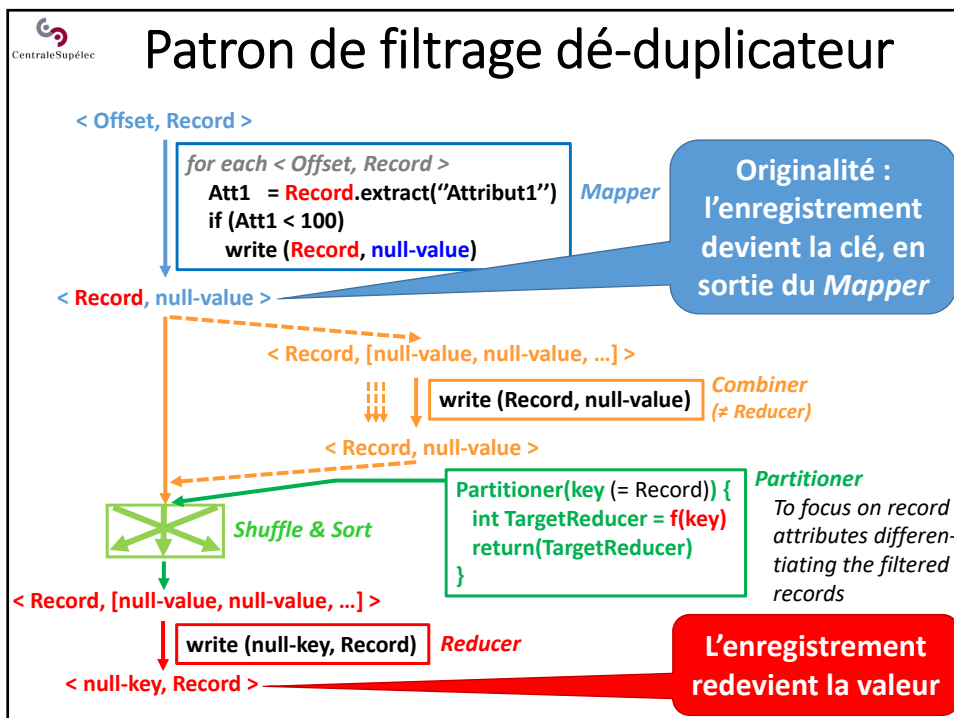
- 1 – Patron de filtrage accepter/rejeter
- 2 – Patron de filtrage *Top Ten*
- 3 – Patron de filtrage dé-duplicateur

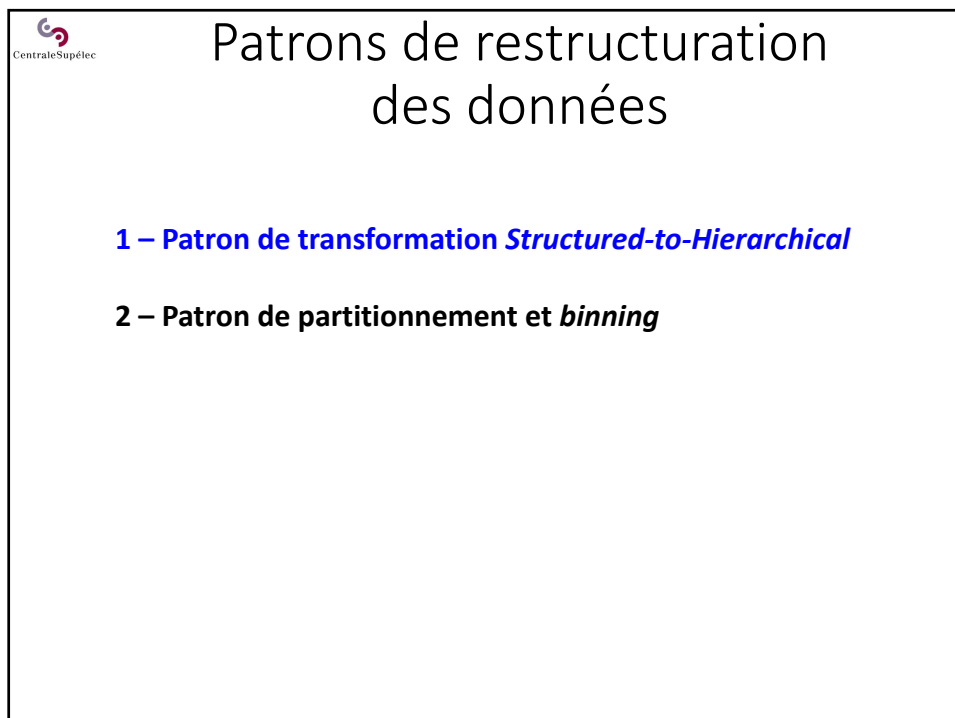
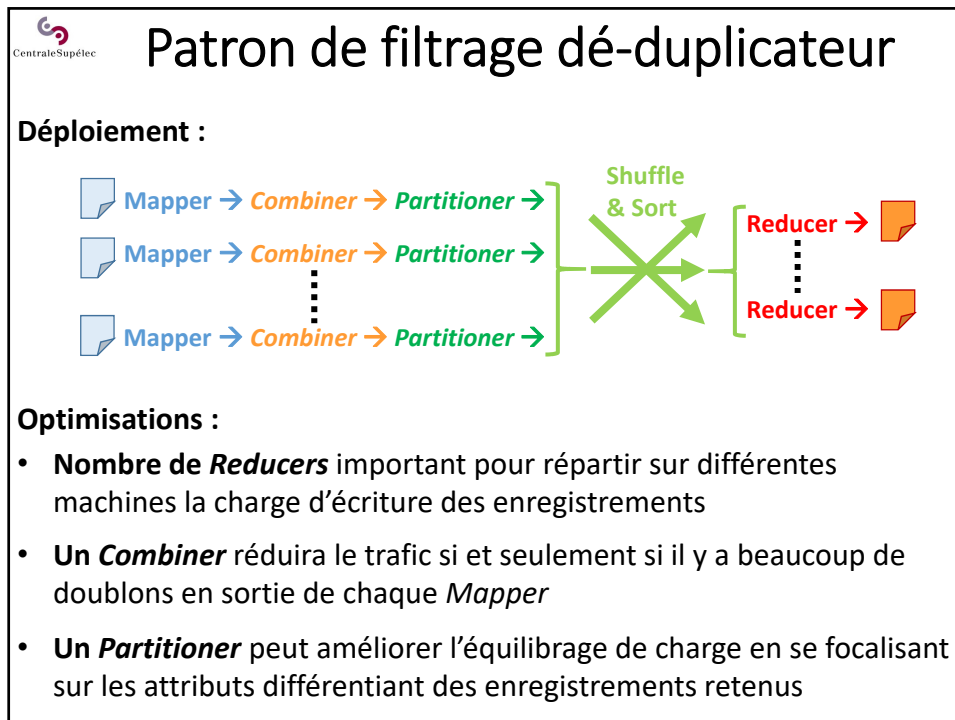


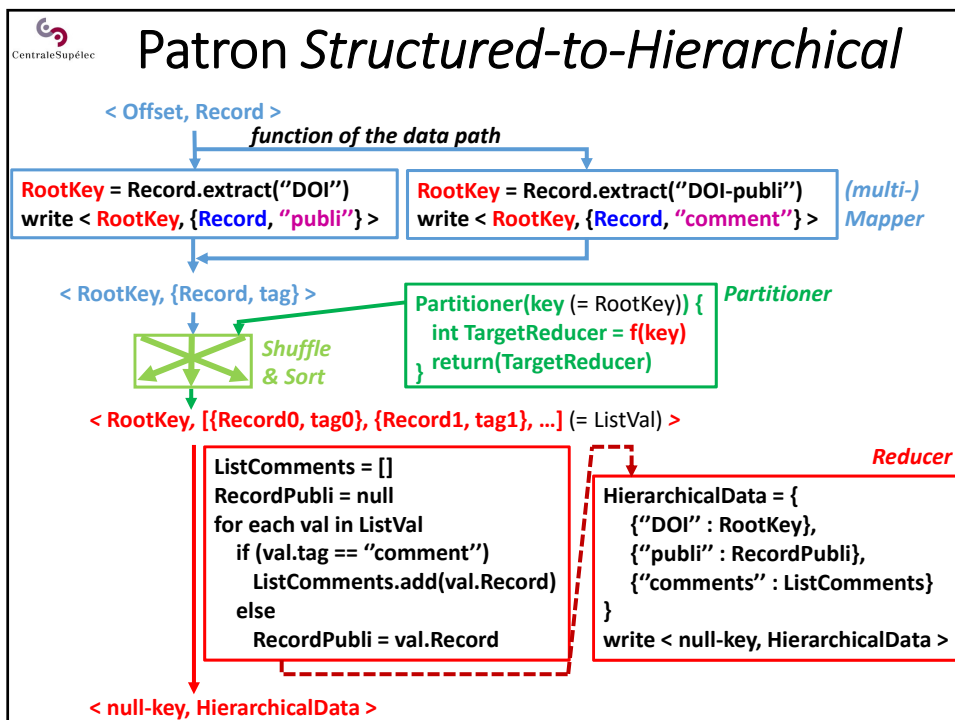
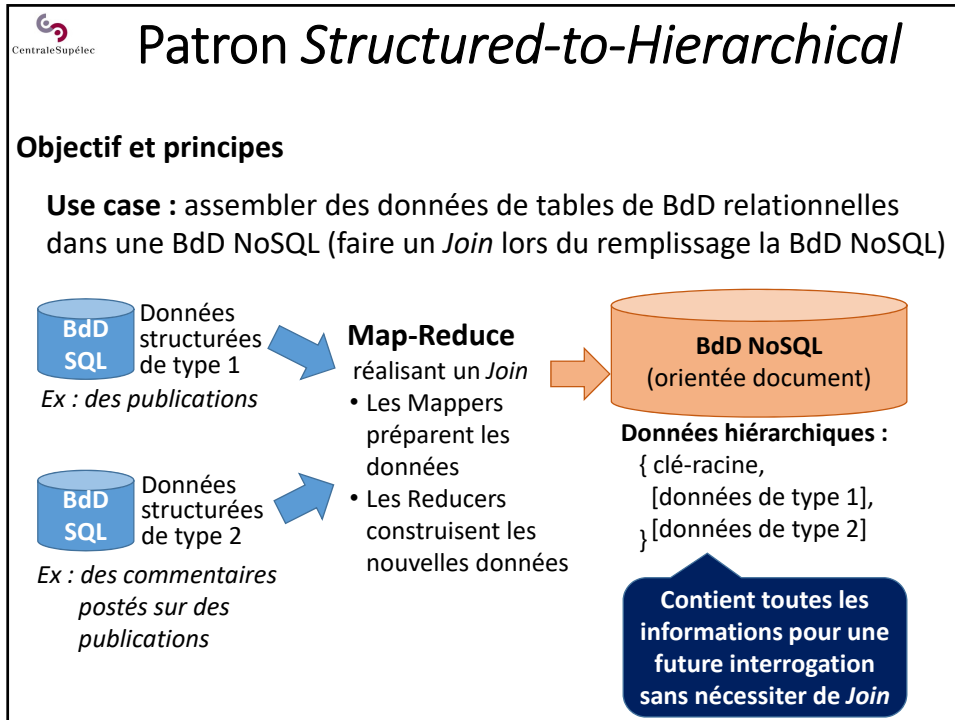
CentraleSupélec

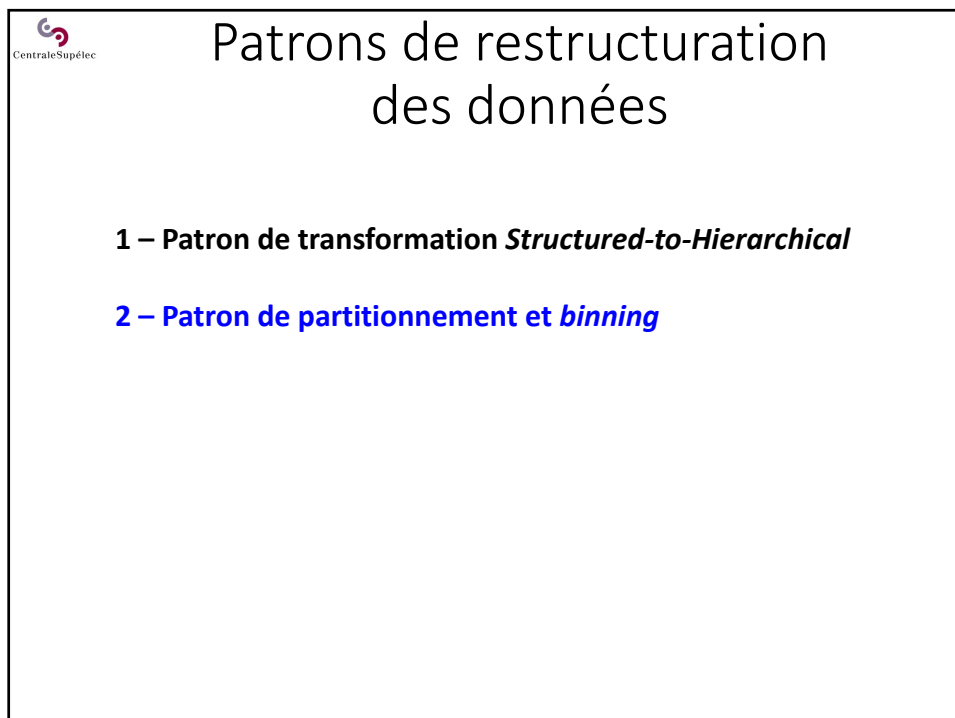
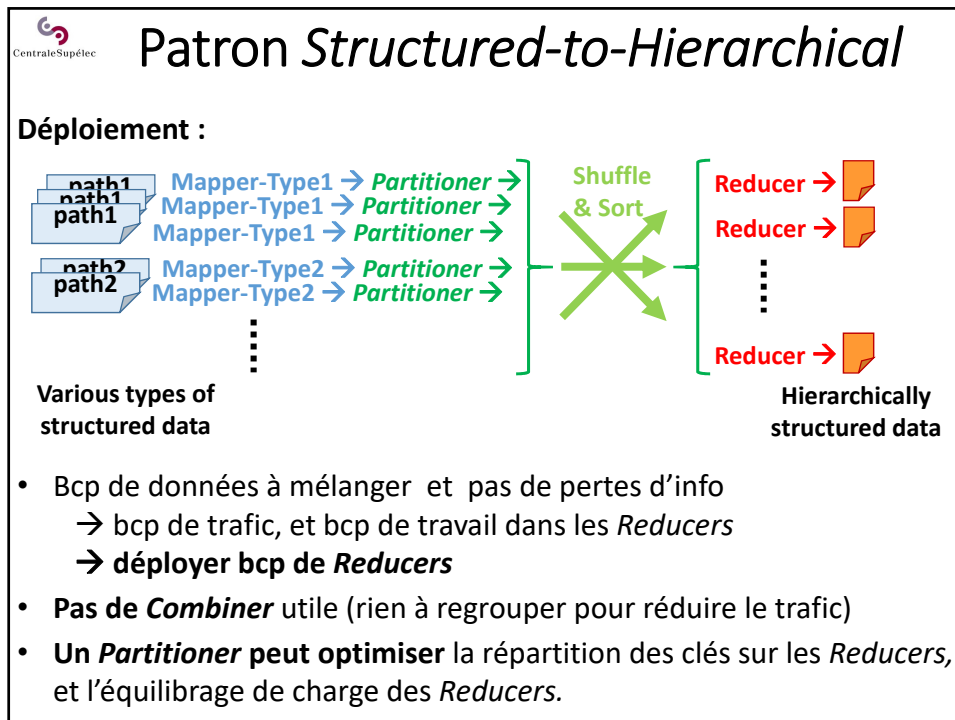
Patrons de Filtrage

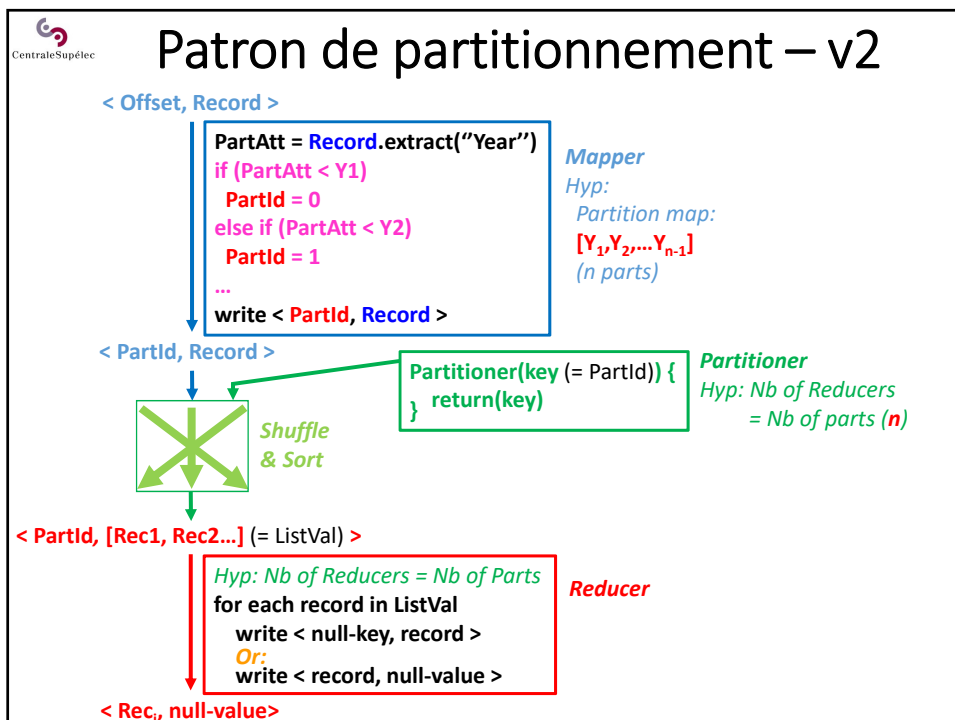
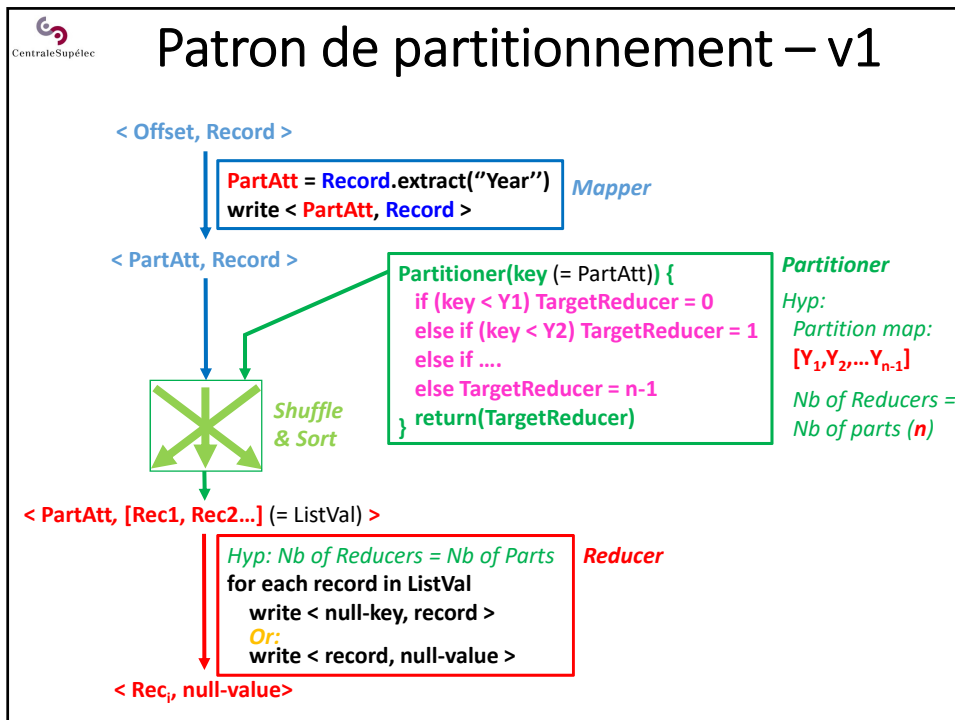
- 1 – Patron de filtrage accepter/rejeter
- 2 – Patron de filtrage *Top Ten*
- 3 – Patron de filtrage dé-duplicateur











CentraleSupélec

Patron de partitionnement

Déploiement :

- Pas de *Combiner* utile (rien à regrouper pour réduire le trafic)
- Toute donnée lu doit être ré-écrite → bcp d'écritures sur disque
- Démarche la plus simple :
 - Un *Reducer* par partition, pour sauver chaque partition dans un fichier différent
 - Un *Partitioner* qui envoie la partition *i* sur le *Reducer i*

CentraleSupélec

Patron de partitionnement

Variante du Patron de Binning (ou de mise en récipients) :

Une donnée peut aller dans plusieurs partitions

Ex : les enregistrements à la limite de deux sous-ensembles sont rangés dans les deux !

Adaptation du *Mapper* de la 2^{ème} version du patron de partitionnement :

```

PartAtt = Record.extract("Attribute1")
if (PartAtt < V0)
  write < 0, Record >
else if (PartAtt == V0)
  write < 0, Record >
  write < 1, Record >
else if (PartAtt < V1)
  write < 1, Record >
...

```

Mapper

Les Mappers peuvent générer plusieurs paires de sorties pour un même enregistrement d'entrée

CentraleSupélec

Patron de partitionnement

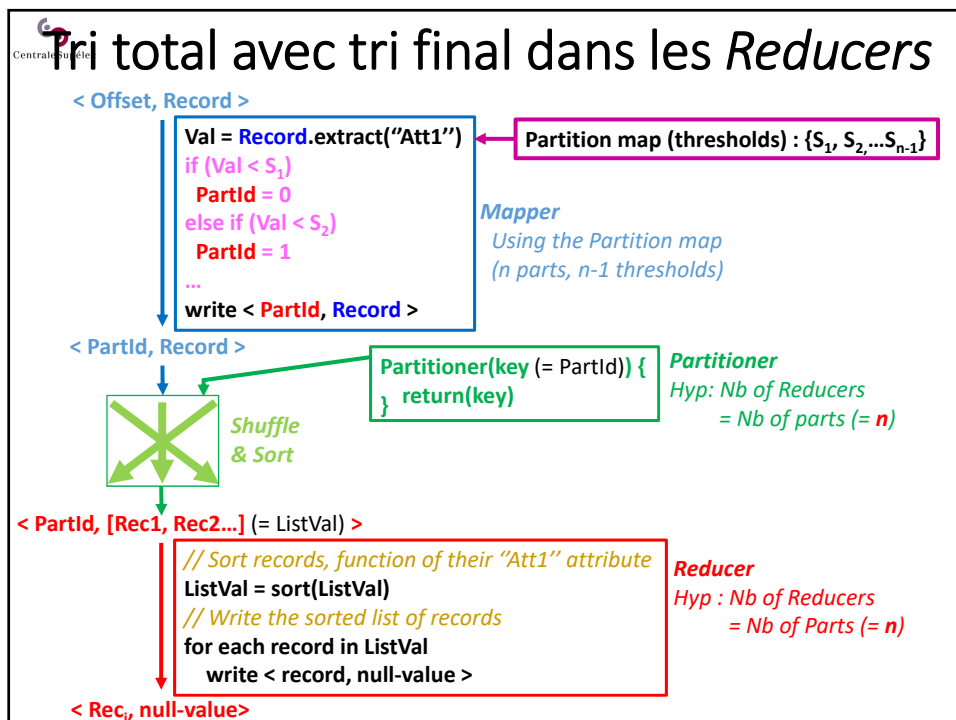
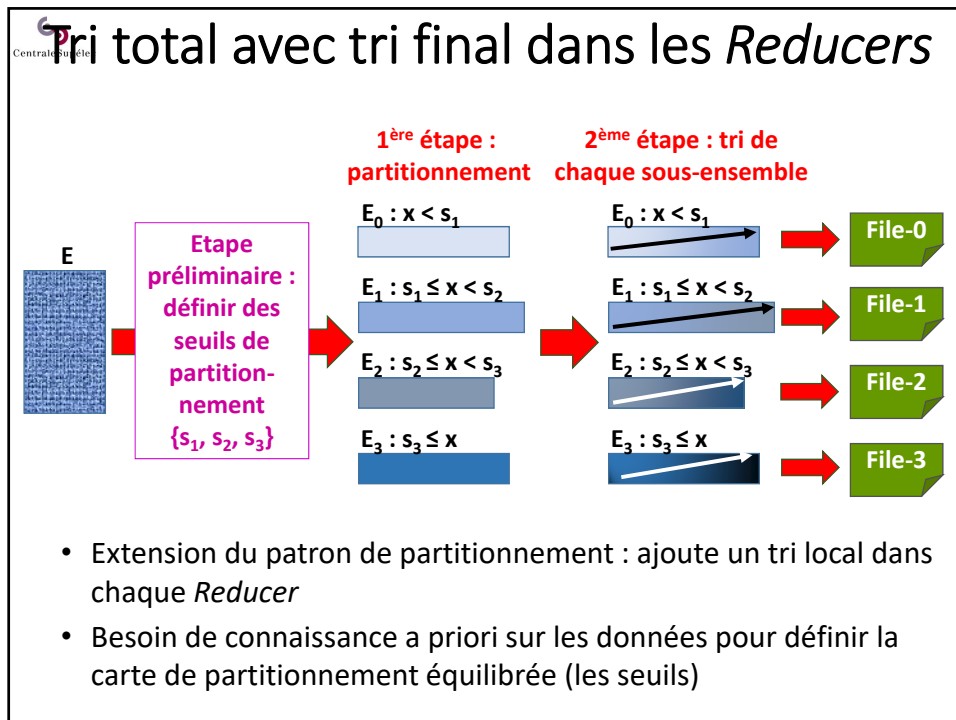
Déploiement :

- Pas de *Combiner* utile (rien à regrouper pour réduire le trafic)
- Toute donnée lu doit être ré-écrite → bcp d'écritures sur disque
- Autre démarche :
 - plus de *Reducers* que de partitions, et ...
 - ...le *Partitioner* répartit les grosses partitions sur plusieurs *Reducer*

CentraleSupélec

Patrons de tri de données

- 1 - Tri total avec tri final dans les *Reducers*
- 2 - Reconfiguration du *Shuffle & Sort* pour optimisation
- 3 - Tri optimisé par le *Shuffle & Sort* sur clés simples et continues
- 4 - Tri optimisé par le *Shuffle & Sort* sur clés composites



CentraleSupélec

Tri total avec tri final dans les *Reducers*

Déploiement et bilan

Partition map (thresholds) : $\{S_1, S_2, \dots, S_{n-1}\}$

Mapper → Partitioner → Shuffle & Sort → Reducer → 1 *Reducer* per subset in the partition (basic deployment)

- Suit le principe algorithmique de base du tri total (1+2 étapes)
- Nbr de *Reducers* = Nbr de sous-ensembles de la partition
- **Ne tire pas partie des tris systématiques et implicites sur les clés à l'entrée des *Reducers*...**

CentraleSupélec

Patrons de tri de données

- 1 - Tri total avec tri final dans les *Reducers*
- 2 - Reconfiguration du *Shuffle & Sort* pour optimisation
- 3 - Tri optimisé par le *Shuffle & Sort* sur clés simples et continues
- 4 - Tri optimisé par le *Shuffle & Sort* sur clés composites

CentraleSupélec

Reconfiguration du *Shuffle & Sort*

Concept de clé composite

Clé simple : $\langle k, v \rangle$
 k : une seule valeur (ex : 10, "toto"...)

Clé composite : ex : $\langle \{k1, k2, k3\}, v \rangle$
 clé composite de 3 attributs

k1 : une valeur (ex : 10, "toto"...)
 k2 : une autre valeur (ex : 'A')
 k3 : une autre valeur (ex : "hf5676klkhjujgh798")

En général, les attributs d'une clé composite correspondent à des critères de tris emboîtés des enregistrements

Ex : trier selon k1 en ordre croissant, puis les enregistrements selon k2 en ordre décroissant, puis ...

CentraleSupélec

Reconfiguration du *Shuffle & Sort*

1^{er} niveau de contrôle/reconfiguration : redéfinition du *Partitioner*

Default *Partitioner*:
 $\text{Reducer Id} = \text{hcode}(k) \% \text{NbReducers}$

CentraleSupélec

Reconfiguration du *Shuffle & Sort*

1^{er} niveau de reconfiguration : **redéfinition du *Partitioner***

```
public static class MyPartitioner
  extends Partitioner<keyClass,valueClass> {
  ...
  public int getPartition(keyClass key,
                        valueClass val,
                        int numPartitions) {
    ...
    // Must return a Reducer index [0 ... NbReducers-1]
    int index = ... // Specific code
    return(index)
  }
}
```

```
//Set the new partitioner to the Map-Reduce job
job.setPartitionerClass{MyPartitioner.class}

// Run the Map-Reduce job
System.exit(job.waitForCompletion(true) ? 0 : 1);
```

CentraleSupélec

Reconfiguration du *Shuffle & Sort*

2^{ème} niveau de reconfiguration : **redéfinition du *keyComparator***
(ou *sortComparator*)

Hyp : clés simples
 $k1 < k2 < k3 < k4 < k5$

Un mapper → $\langle k3, v3 \rangle$

Un mapper → $\langle k2, v2 \rangle$

Un mapper → $\langle k4, v4 \rangle$

keyComparator → Reducer i

Sorted key-value pairs:

- $\langle k1, v1 \rangle$
- $\langle k2, v2 \rangle$
- $\langle k3, v3 \rangle$
- $\langle k4, v4 \rangle$
- $\langle k5, v5 \rangle$

CentraleSupélec

Reconfiguration du *Shuffle & Sort*

2^{ème} niveau de reconfiguration : **redéfinition du *keyComparator***
(ou *sortComparator*)

Hyp : clés composites
 $k_1 < k_2 < k_3$
 $t_{11} < t_{12}$, et $t_{21} < t_{22}$

Un mapper → $\langle \{k_2, t_{21}\}, v_3 \rangle$
 Un mapper → $\langle \{k_1, t_{11}\}, v_1 \rangle$
 Un mapper → $\langle \{k_1, t_{12}\}, v_2 \rangle$ & $\langle \{k_3, t_{31}\}, v_5 \rangle$
 Un mapper → $\langle \{k_2, t_{22}\}, v_4 \rangle$

keyComparator Reducer i

keyComparator sur les 2 attributs de la clé composite

$\langle \{k_1, t_{11}\}, v_1 \rangle$
 $\langle \{k_1, t_{12}\}, v_2 \rangle$
 $\langle \{k_2, t_{21}\}, v_3 \rangle$
 $\langle \{k_2, t_{22}\}, v_4 \rangle$
 $\langle \{k_3, t_{31}\}, v_5 \rangle$

CentraleSupélec

Reconfiguration du *Shuffle & Sort*

2^{ème} niveau de reconfiguration : **redéfinition du *keyComparator***
(ou *sortComparator*)

```

public static class MyKeyComparator
  extends WritableComparator {
  ...
  public int compare(WritableComparable w1,
                    WritableComparable w2) {
    ...
    // Must return -1/0/+1,
    // considering w1 < w2, or w1 = w2, or w1 > w2
    int ComparisonResult = ... // Specific code
    return (ComparisonResult)
  }
}

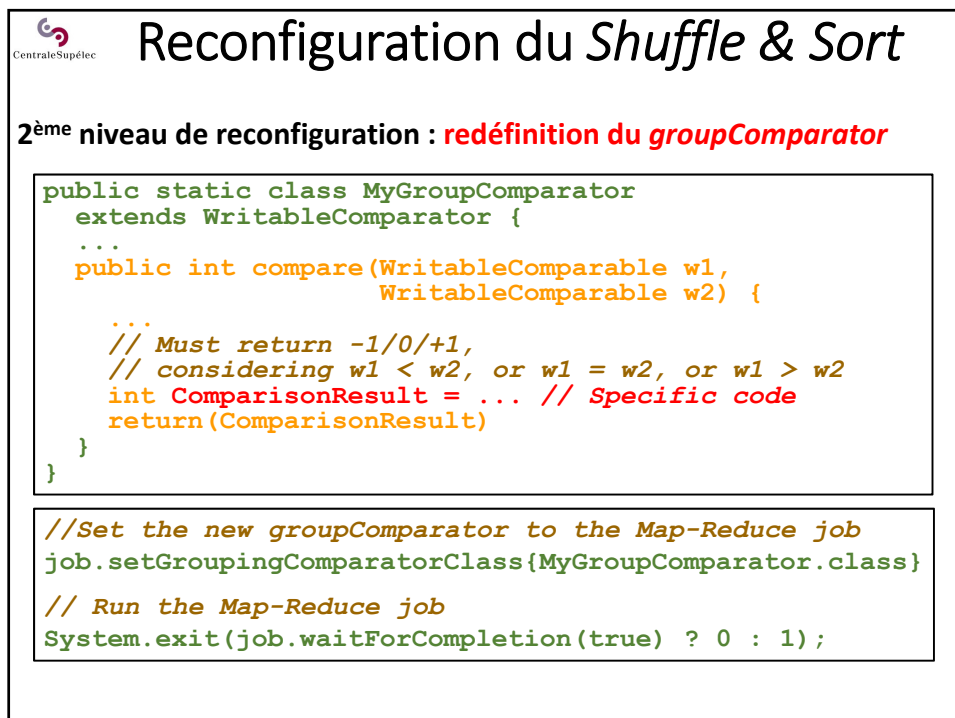
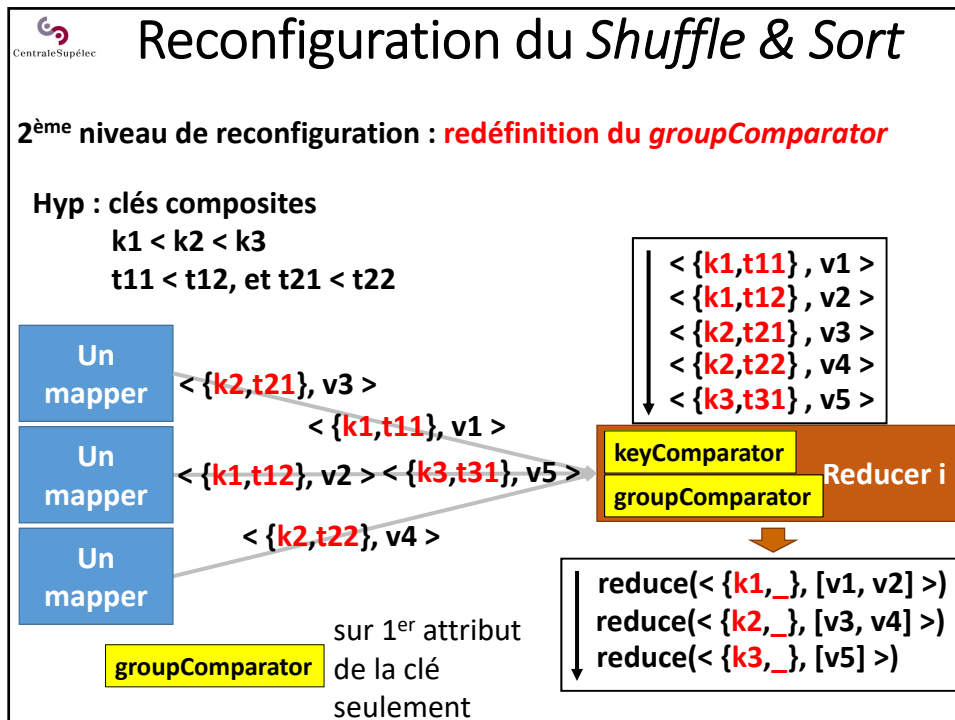
```

```

//Set the new keyComparator to the Map-Reduce job
job.setSortComparatorClass(MyKeyComparator.class)

// Run the Map-Reduce job
System.exit(job.waitForCompletion(true) ? 0 : 1);

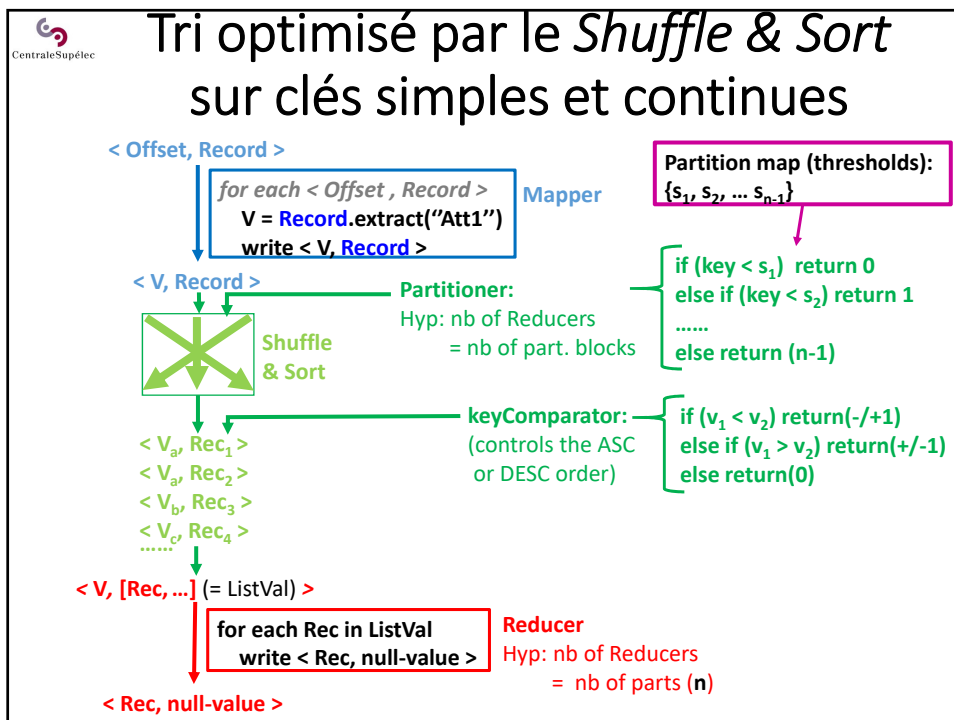
```



CentraleSupélec

Patrons de tri de données

- 1 - Tri total avec tri final dans les *Reducers*
- 2 - Reconfiguration du *Shuffle & Sort* pour optimisation
- 3 - Tri optimisé par le *Shuffle & Sort* sur clés simples et continues
- 4 - Tri optimisé par le *Shuffle & Sort* sur clés composites



CentraleSupélec

Tri optimisé par le *Shuffle & Sort* sur clés simples et continues

Déploiement et bilan

Partition map (thresholds) : $\{S_1, S_2, \dots, S_{n-1}\}$

1 *Reducer* per subset in the partition

- Profite du tri des clés du *Shuffle & Sort*, et évite un tri final dans les *Reducers*
- Nbr de *Reducers* = Nbr de sous-ensembles de la partition
- Ne trie les enregistrements que selon UN SEUL critère

CentraleSupélec

Patrons de tri de données

- 1 - Tri total avec tri final dans les *Reducers*
- 2 - Reconfiguration du *Shuffle & Sort* pour optimisation
- 3 - Tri optimisé par le *Shuffle & Sort* sur clés simples et continues
- 4 - Tri optimisé par le *Shuffle & Sort* sur clés composites

CentraleSupélec

Tri optimisé par le *Shuffle & Sort* sur clés composites

Problématique / Objectifs

- **Trier selon deux critères (ou plus) emboîtés :**
SELECT * FROM etudiant ORDER BY spécialité, moyenne DESC
- **Faire un traitement sur les sous-listes du dernier critère**
Ex : calcul d'une valeur moyenne médiane ou maximale par spécialité
- **Profiter du tri implicite et systématique d'Hadoop sur les clés**

Démarche

- Construire une **clé composite** avec les deux critères :
Ex : {spécialité, moyenne}
- Redéfinir les **3 niveaux de contrôles** du *Shuffle & Sort*

CentraleSupélec

Tri optimisé par le *Shuffle & Sort* sur clés composites

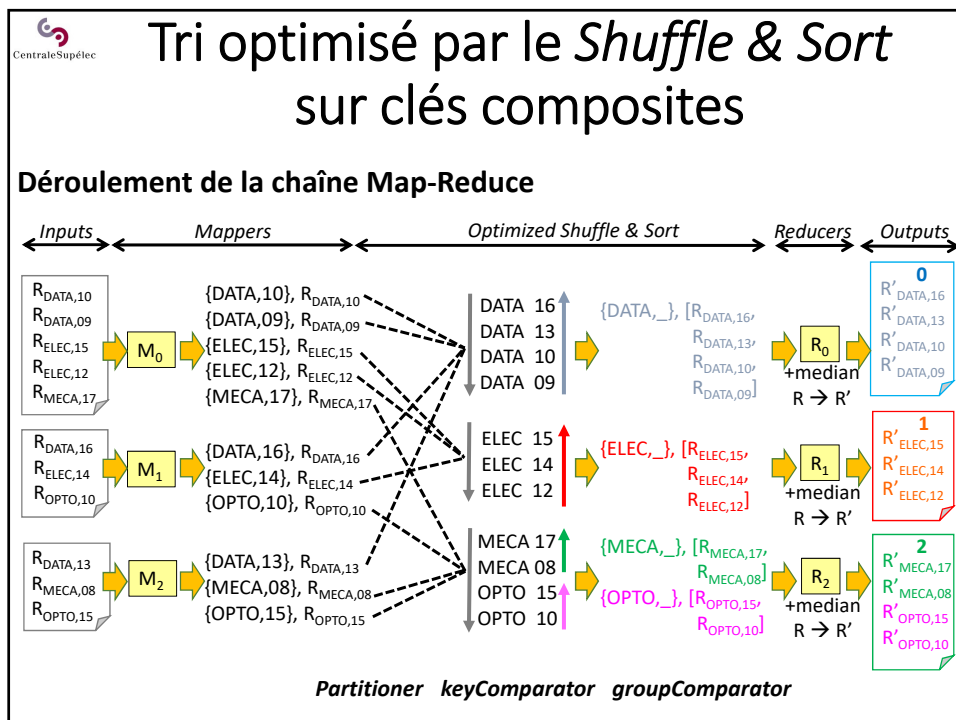
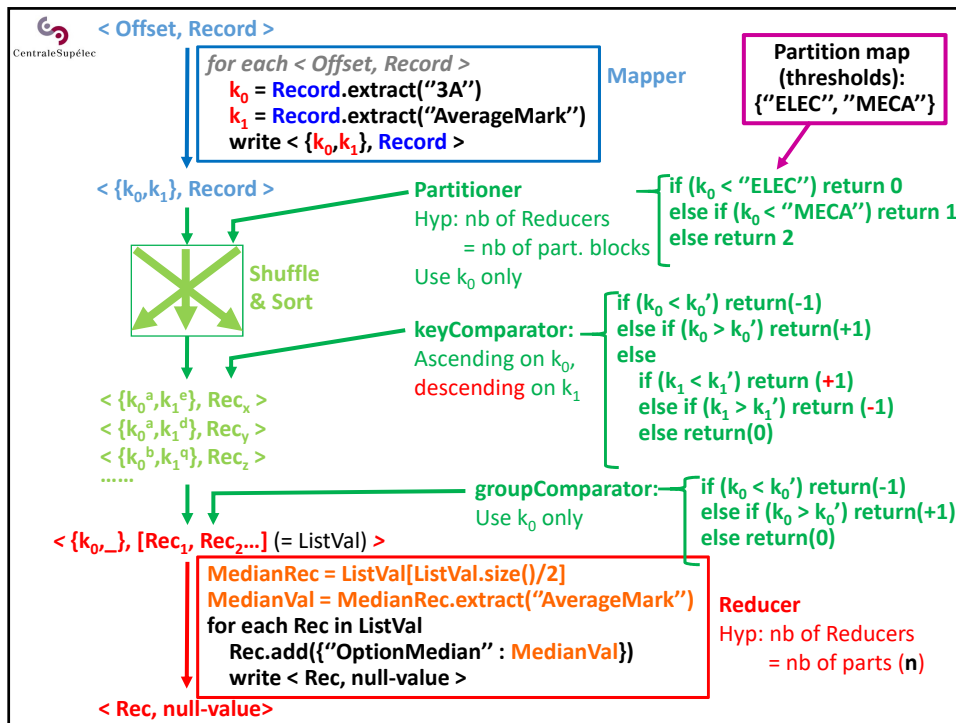
Problème illustratif

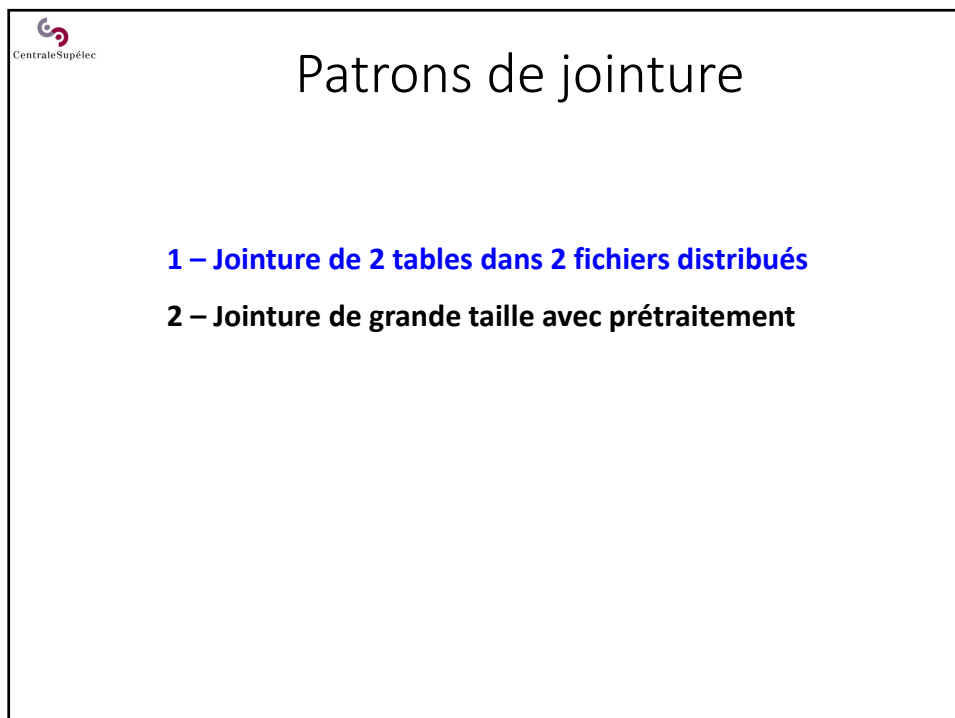
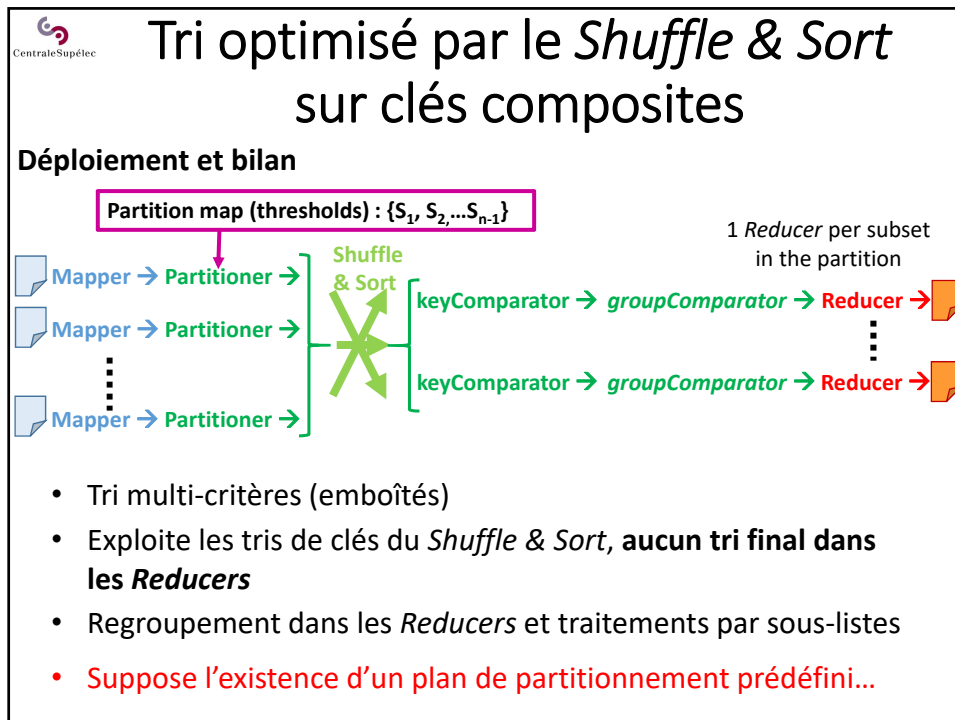
Trier des enregistrements d'étudiants **selon 2 critères** emboîtés :

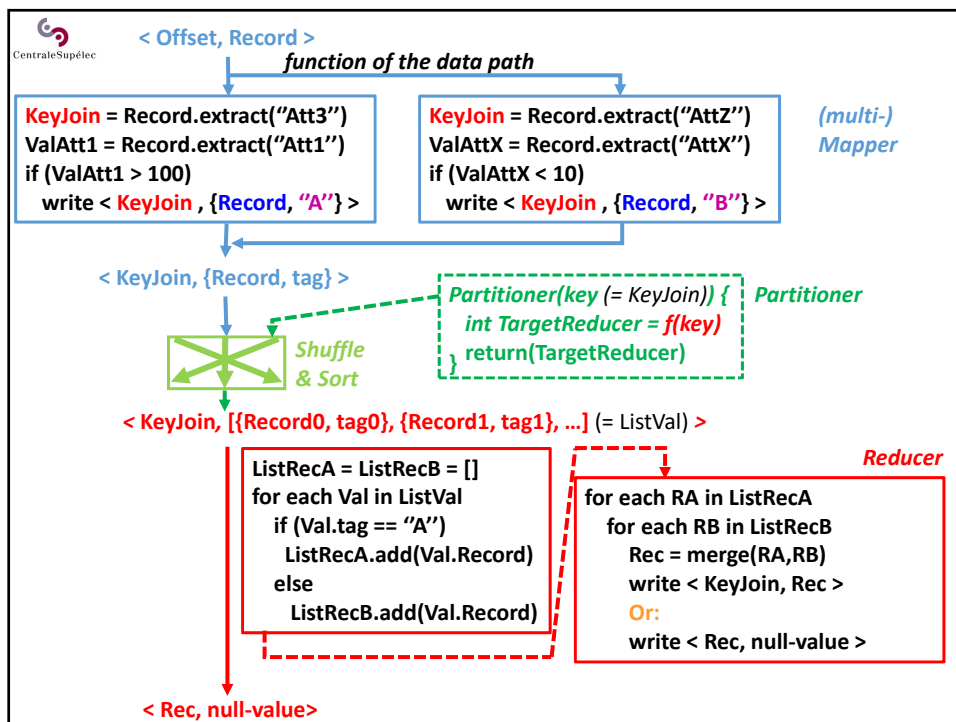
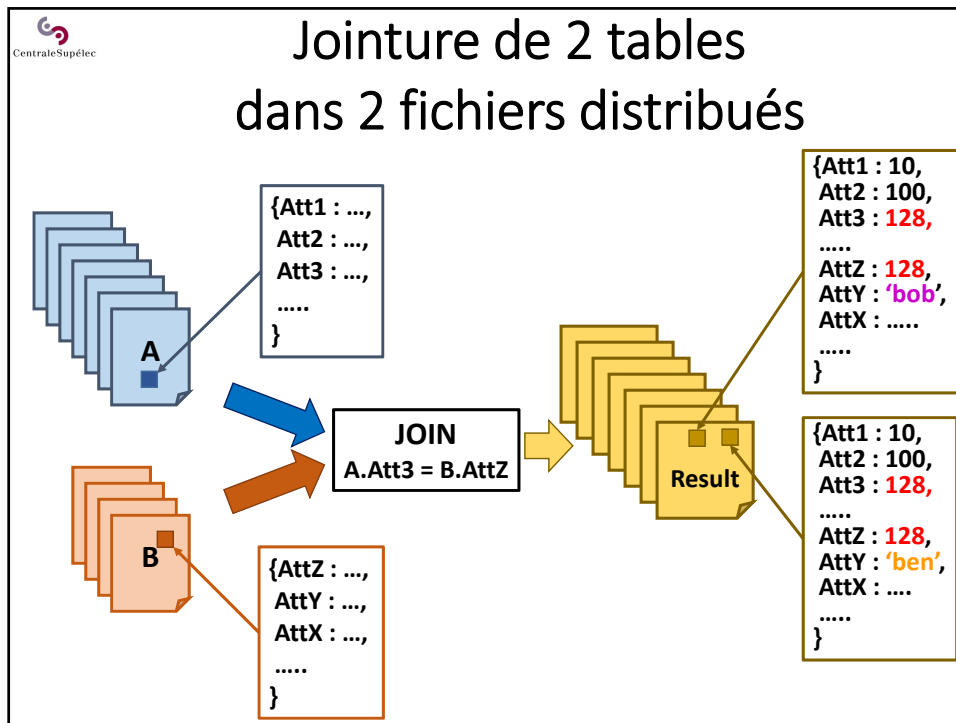
- 1 – par option de 3A en ordre lexicographique croissant
- 2 – par moyenne générale en ordre décroissant

Pour chaque option de 3A :

- calculer la médiane des moyennes générales
- enrichir chaque enregistrement avec cette valeur médiane
- écrire sur disque dans l'ordre les enregistrements enrichis







CentraleSupélec

Patrons de jointure


- 1 – Jointure de 2 tables dans 2 fichiers distribués
- 2 – Jointure de grande taille avec prétraitement

CentraleSupélec

Jointure de grande taille avec prétraitement

Principe

1. Identifier un format de donnée structurée à partir duquel il serait simple et rapide de faire un *A JOIN B*
Supposer que les attributs de l'équi-jointure sont connus et fixés
2. Implanter le pré-traitement (en *Map-Reduce*) pour obtenir de nouvelles versions des documents structurés/tables *A* et *B*
→ réaliser un « *co-partitionnement* »...
3. Implanter l'opération *JOIN* (en *Map-Reduce*) sur les nouvelles versions des tables *A* et *B*

 Exercice



Outils d'aide au partitionnement

- 1 – *TotalOrderPartitioner*
- 2 – Configuration d'un partitionneur par échantillonnage
- 3 – Echantillonneurs prédéfinis



Classe *TotalOrderPartitioner*

La classe Java d'Hadoop *TotalOrderPartitioner* implémente un *Partitioner* qui exploite une **carte de partitionnement prédéfinie**

```
// Create the 'Map-Reduce job'
Job MRJob = new Job(getConf());
// Configure the 'Map-Reduce job'
// - Set the Mapper and Reducer to use
MRJob.setMapperClass(MyMapper.class)
MRJob.setReducerClass(Myreducer.class)

// - Set the number of Reducers
MRJob.setNumReduceTasks(REDUCE_TASKS);
// - Set the Partitioner to use: 'TotalOrderPartitioner'
MRJob.setPartitionerClass(TotalOrderPartitioner.class);
// - Set the partition map used by the 'TotalOrderPartitioner' (set partition file path)
TotalOrderPartitioner.setPartitionFile(MRJob.getConfiguration(), partitionPath);

// Run the Map-Reduce job, using the TotalOrderPartitioner
System.exit(MRJob.waitForCompletion(true) ? 0 : 1);
```

→ Il faut disposer d'une **carte de partitionnement prédéfinie** et menant à un bon équilibrage de charge !!



Outils d'aide au partitionnement

1 – *TotalOrderPartitioner*

2 – Configuration d'un partitionneur par échantillonnage

3 – Echantillonneurs prédéfinis



Configuration d'un partitionneur par échantillonnage

Principe

Un **partitionnement équilibré** évite qu'UN REDUCER reçoive et traite beaucoup plus de valeurs que les autres :

- serait couteux en temps de traitement
- risquerait de saturer la mémoire du *Reducer*

Mais sans connaissance a priori sur les données ni expérience sur des données similaires...on ne sait pas établir une carte de partitionnement équilibrée !

→ **Il faut analyser le jeu de données**
Long ! Et demande beaucoup de moyens !

→ **Analyser un échantillon du jeu de données**

CentraleSupélec

Configuration d'un partitionneur par échantillonnage

Principe

→ Analyser un échantillon du jeu de données

- Définir une politique d'échantillonnage
- Implanter un échantillonneur qui accède aux données sur HDFS

Hadoop fournit des outils d'échantillonnage de données et de construction d'une carte de partitionnement équilibrée

Exemples :

- Echantillonnage aléatoire
- Echantillonnage à intervalles réguliers
- Collecte des n premières valeurs

CentraleSupélec

Configuration d'un partitionneur par échantillonnage

Principe

→ Analyser un échantillon du jeu de données

- Définir une politique d'échantillonnage
- Implanter un échantillonneur qui accède aux données sur HDFS

Hadoop fournit des outils d'échantillonnage de données et de construction d'une carte de partitionnement équilibrée

+ simples à utiliser

- Tournent dans l'appli cliente sur une seule machine (la solution sort du paradigme Map-Reduce)



Outils d'aide au partitionnement

- 1 – *TotalOrderPartitioner*
- 2 – Configuration d'un partitionneur par échantillonnage
- 3 – **Echantillonneurs prédéfinis**



Echantillonneurs prédéfinis (1)

```
// Define and start to configure a new Map-Reduce job
Job job = ...
job.setXXXX...
job.setYYYY...

// - Set the number of Reducers (required before to run the sampler)
job.setNumReduceTasks(REDUCE_TASKS);
// - Set the partition output path to use for the Map-Reduce job
// (do it across a Partitioner configuration)
TotalOrderPartitioner.setPartitionFile(job.getConfiguration(), partitionOutputPath);

// Choice and configuration of a RandomSampler
// 10% of chance to choose a data, 10000: max number of choosen data
// 10: max number of split analyzed
InputSampler.Sampler mySmpl = new InputSampler.RandomSampler(0.1,10000,10);

// Run the Sampler that builds and writes the partition map:
// - job is the Map-Reduce job requiring the partition map
// - Its reference allows to get the number of planned Reducers, the input data path,
// the partition output path...
InputSampler.WritePartitionFile(job, mySmpl);
```

CentraleSupélec

Echantillonneurs prédéfinis (2)

```

TotalOrderPartitioner.setPartitionFile(job.getConfiguration(), partitionOutputPath);
InputSampler.Sampler mySmpl = new InputSampler.RandomSampler(0.1,10000,10);
InputSampler.WritePartitionFile(job, mySmpl);
// Complete the Map-Reduce job configuration
job.setPartitionerClass(TotalOrderPartitioner.class);
job.setMapperClass(MyMapper.class);
job.setReducerClass(MyReducer.class);
...
// Run the Map-Reduce job, using the Partitioner that uses the partition map
// computed by the sampler.
System.exit(job.waitForCompletion(true) ? 0 : 1);

```

```

graph TD
    A[Definir : un job Map-Reduce  
un partitioner  
un échantillonneur] --> B[Les configurer et les associer]
    B --> C[Exécuter l'échantillonneur qui génère  
la carte de partitionnement]
    C --> D[Exécuter le job Map-Reduce avec le  
partitionneur qui exploite la carte]

```

CentraleSupélec

Bilan



Bilan du Map-Reduce d'Hadoop

Finalemment :

- Les *Mappers* peuvent sortir du cadre du SPMD (multi-mappers) quand on utilise plusieurs fichiers (test sur le *path* des fichiers)
→ Gagne en généralité
- Le *Shuffle & Sort* semble être un schéma de comm contraint. Mais c'est une sorte de all-to-all des *Mappers* vers les *Reducers*, et reconfigurable en 3 points !
→ C'est un schéma de comm assez générale
- Un pipelining recouvre les calculs des *Mappers*, les communications et les regroupements de données
→ Le *Map-Reduce* possède une implantation optimisée
- Le paradigme *Map-Reduce* masque les aspects d'informatique distribuée aux développeurs applicatifs
→ Permet l'adhésion d'une large "communauté d'utilisateurs Big Data" (à l'opposé de la stratégie du HPC)



Algorithmique Map-Reduce

