





CentraleSupélec  

Big Data : Informatique pour les données et calculs massifs

## 6 – Algorithmique Map-Reduce

**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

## Démarche algorithmique

- **Définir des tâches *Mappers* et *Reducers***  
→ Bâtir la base de la solution *Map-Reduce*
- **Identifier des routines *Combiner***  
→ Eviter des E/S temporaires sur disque (*Hadoop*)  
→ Réduire le trafic sur le réseau entre *Mappers* et *Reducers*  
→ Eviter de saturer la mémoire des *Reducers*
- **Identifier le nombre de *Reducers* optimal**  
→ Répartir la charge au maximum entre plusieurs machines
- **Concevoir des *Partitioners* pour optimiser la solution**  
→ Améliorer la répartition de charge des *Reducers*  
(si on connaît la variété des clés et la distribution des paires)
- **Ecrire des fonctions *keyComparator()* et *groupComparator()***  
→ Pour profiter des fonctions de tri intégrées au *Shuffle & Sort*

➡ « Patterns » ou « Patrons » de conception types

---

---

---

---

---

---

---

---

CentraleSupélec

## Familles de patrons Map-Reduce

- Patrons de récapitulation (*Summarization*)
- Patrons de filtrage
- Patrons de restructuration des données
- Patrons de tris
- Patrons de jointure

---

---

---

---

---

---

---

---

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

---

---

---

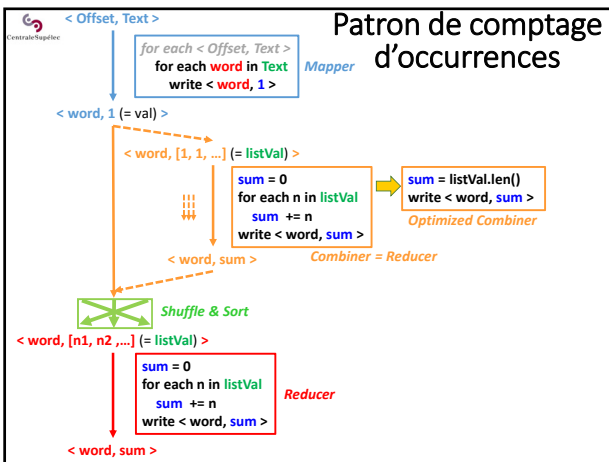
---

---

---

---

---




---

---

---

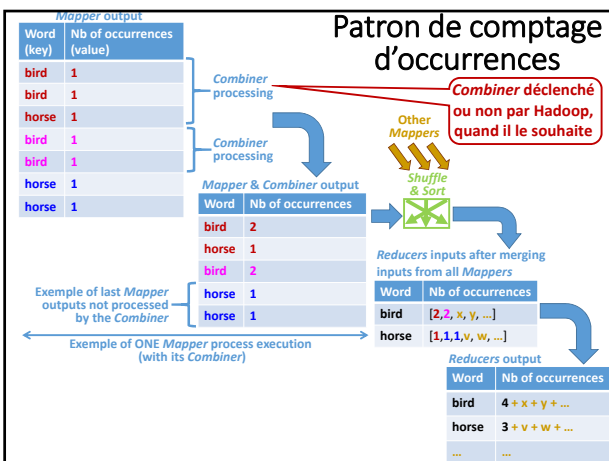
---

---

---

---

---




---

---

---

---

---

---

---

---

### Patron de comptage d'occurrences

**Déploiement :**

- 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!

---

---

---

---

---

---

---

---

---

---

### 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

---

---

---

---

---

---

---

---

---

---

### Patron de calcul d'une moyenne

```

< Offset, Record >
for each < Offset, Record >
  id = Record.extract("ProductId")
  price = Record.extract("price")
  write < id, {price, 1} >
  
```

*Mapper*

```

< id, {SumPrice, 1} (= val) >
  SumPrice = NbPrice = 0
  for each val in listVal
    SumPrice += val.SumPrice
    NbPrice++
  write < id, {SumPrice, NbPrice} >
  
```

*Combiner (≠ Reducer)*

```

< id, {SumPrice, NbPrice} >
  
```

*Shuffle & Sort*

```

< id, [{SumPrice1,NbPrice1},{SumPrice2,NbPrice2},...] (= listVal) >
  Sum = N = 0
  for each val in listVal
    Sum += val.SumPrice
    N += val.NbPrice
  write < id, Sum/N >
  
```

*Reducer*

```

< id, average >
  
```

---

---

---

---

---

---

---

---

---

---

**Patron de comptage d'occurrences**

Déploiement :

Même déploiement que pour le compteur d'occurrences

- 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 produits attendus
- Un *Partitioner* pour améliorer l'équilibrage de charge entre *Reducers* ... ssi connaissances pertinentes sur les résultats attendus!

---

---

---

---

---

---

---

---

---

---

**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

---

---

---

---

---

---

---

---

---

---

**Patron d'index inversé**

```

< Offset, Text >
for each < Offset, Text >
  DocId = Text.extract("DocId")
  for each Word in Text
    if (IsInteresting(Word))
      EntryWord = CanonicForm(Word)
      write < EntryWord, DocId >
  
```

*Mapper*

```

Partitioner(key (= Word)) {
  int TargetReducer = f(key)
  return(TargetReducer)
}
  
```

*Partitioner*  
Only if *knowledge* about word distribution allows to improve the default partitioner function

Shuffle & Sort

```

< Word, [DocId1, DocId2, ...] >
  
```

Reducer

```

Close to identity fct
(output format improvement)
  
```

```

< Word, [DocId1, DocId2, ...] >
  
```

---

---

---

---

---

---

---

---

---

---

## Patron d'index inversé

**Déploiement :**

```

    graph LR
      subgraph Mappers
        M1[Mapper] --> P1[Partitioner]
        M2[Mapper] --> P2[Partitioner]
        M3[Mapper] --> P3[Partitioner]
      end
      P1 --> SS[Shuffle & Sort]
      P2 --> SS
      P3 --> SS
      SS --> R1[Reducer - Identity]
      SS --> R2[Reducer - Identity]
  
```

- **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

---

---

---

---

---

---

---

---

---

---

## Patrons de Filtrage

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

---

---

---

---

---

---

---

---

---

---

## Patron de filtrage accepter/rejeter

```

< Offset, Record >
  for each < Offset, Record >
    RecKey = GenerateUniqueKey()
    Att1 = Record.extract("Attribut1")
    Att2 = Record.extract("Attribut2")
    if (Att1 < 100 and Att2 > 0)
      write (RecKey, Record)
  
```

*Mapper*

```

Partitioner(key (= RecKey)) {
  int TargetReducer = f(key)
  return(TargetReducer)
}
  
```

*Partitioner*  
Only if *knowledge* about RecKey distribution allows to improve the default partitioner function

```

  Shuffle & Sort
  
```

```

< RecKey [Record] >
  write (RecKey, Record)
  
```

*Reducer*

**Requête SQL équivalente :**

```

SELECT *
FROM table
WHERE condition
          
```

---

---

---

---

---

---

---

---

---

---

**Patron de filtrage accepter/rejeter**

```

< Offset, Record >
  for each < Offset, Record >
    RecKey = GenerateUniqueKey()
    Att1 = Record.extract("Attribut1")
    Att2 = Record.extract("Attribut2")
    if (Att1 < 100 and Att2 > 0)
      write (RecKey, Record)
  
```

Besoin de générer une clé unique pour chaque enregistrement

```

Partitioner(key = RecKey) {
  int TargetReducer = f(key)
  return(TargetReducer)
}
  
```

Fournir le *Partitioner* qui répartit équitablement les clés sur beaucoup de *Reducers*

→ Générateur de clé et *Partitioner* conçus en cohérence

```

Shuffle & Sort
  
```

```

< RecKey [Record] >
  write (RecKey, Record)
  
```

< RecKey, Record >

---

---

---

---

---

---

---

---

---

---

**Patron de filtrage accepter/rejeter**

Déploiement :

```

Mapper → Partitioner → Shuffle & Sort → Reducer (close to Identity)
Mapper → Partitioner → Shuffle & Sort → Reducer (close to Identity)
Mapper → Partitioner → Shuffle & Sort → Reducer (close to Identity)
  
```

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

---

---

---

---

---

---

---

---

---

---

**Patrons de Filtrage**

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

---

---

---

---

---

---

---

---

---

---

**Patron de filtrage Top Ten / Top K**

**Principe et déploiement :**

```

    graph LR
      subgraph Mappers
        M1[Mapper Filtering, sorting and writing the K best records]
        M2[Mapper Filtering, sorting and writing the K best records]
      end
      Mappers -- Shuffle & Sort --> R[Reducer Sorting n x K records and writing the K best records]
  
```

- Chaque *Mapper* filtre des enregistrements, tri ceux acceptés, et écrit les paires clé-valeur des K meilleurs locaux
- UN *Reducer* reçoit les  $n \times K$  meilleurs, les interclasse, et écrit les paires des K meilleurs globaux
- Pas d'optimisations :
  - pas de *Combiner* possible (sorties de chaque *Mapper* déjà réduites),
  - pas de *Partitioner* possible (un seul *Reducer* cible)

---

---

---

---

---

---

---

---

---

---

**Patron de filtrage Top Ten / Top K**

```

    < Offset, Record >
    List<Record> ListRec
    for each < Offset, Record >
      Att1 = Record.extract("Attribut1")
      if (Att1 < 100)
        ListRec.add(Record)
      Sort(ListRec)
      for (i = 0; i < K and i < ListRec.size; i++)
        write (null-key, ListRec[i])
  
```

*Mapper (not the most efficient code)*

*Shuffle & Sort*

```

    < null-key, [Rec0, Rec1, Rec2, ...] (= ListVal) >
    Sort(ListVal)
    for (i = 0; i < K and i < ListVal.size; i++)
      write (i, ListVal[i])
  
```

*Reducer*

*< rank, Record >*

---

---

---

---

---

---

---

---

---

---

**Patrons de Filtrage**

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

---

---

---

---

---

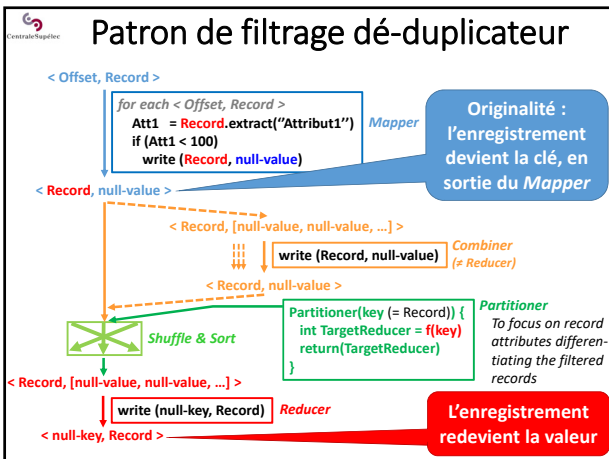
---

---

---

---

---




---

---

---

---

---

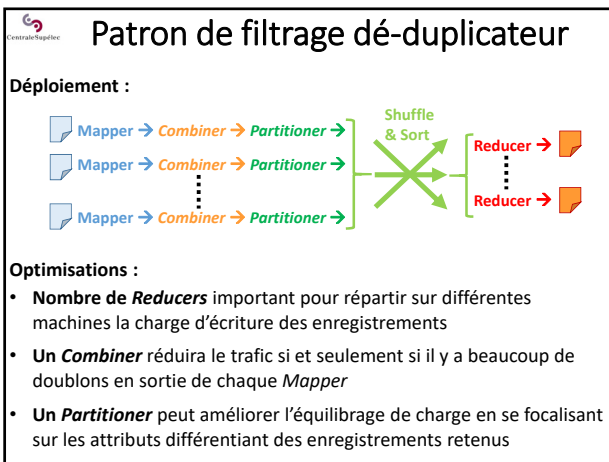
---

---

---

---

---




---

---

---

---

---

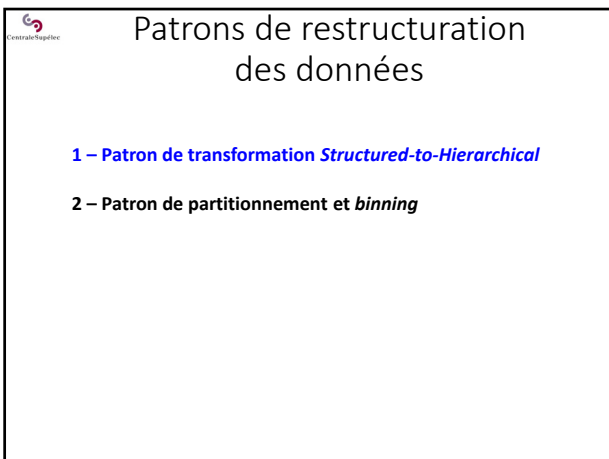
---

---

---

---

---




---

---

---

---

---

---

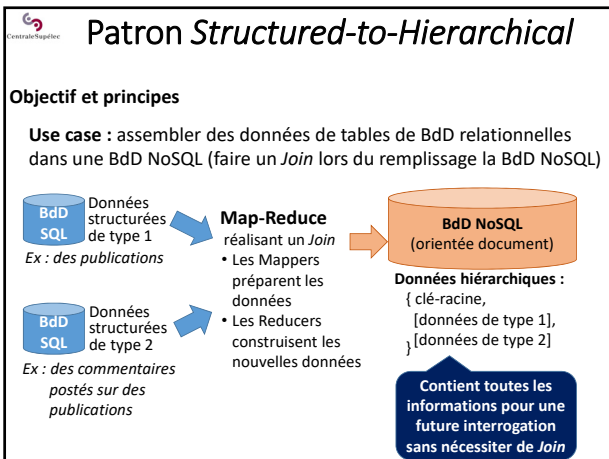
---

---

---

---






---

---

---

---

---

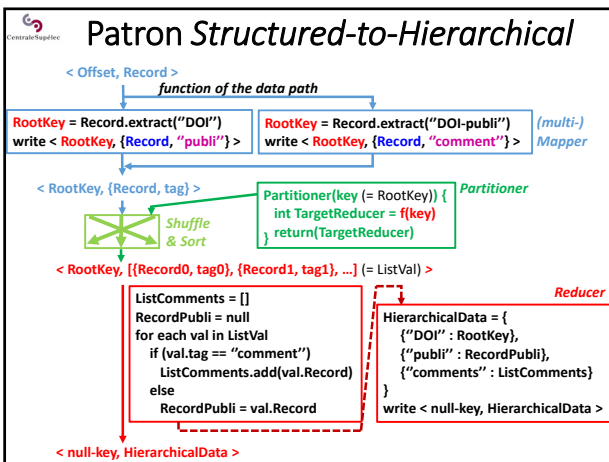
---

---

---

---

---




---

---

---

---

---

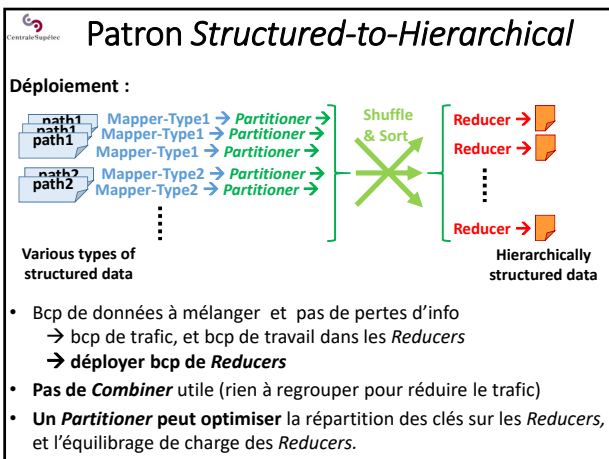
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

CentraleSupélec

## Patrons de restructuration des données

1 – Patron de transformation *Structured-to-Hierarchical*

2 – Patron de partitionnement et *binning*

---

---

---

---

---

---

---

---

CentraleSupélec

## Patron de partitionnement – v1

< Offset, Record >

```

PartAtt = Record.extract("Year")
write < PartAtt, Record >
    
```

Mapper

< PartAtt, Record >

```

Partitioner(key (= PartAtt)) {
  if (key < Y1) TargetReducer = 0
  else if (key < Y2) TargetReducer = 1
  else if ....
  else TargetReducer = n-1
  return(TargetReducer)
}
    
```

Partitioner

Hyp: Partition map: [Y<sub>1</sub>, Y<sub>2</sub>, ..., Y<sub>n-1</sub>]

Nb of Reducers = Nb of parts (n)

Shuffle & Sort

< PartAtt, [Rec1, Rec2...] (= ListVal) >

```

Hyp: Nb of Reducers = Nb of Parts
for each record in ListVal
  write < null-key, record >
  Or:
  write < record, null-value >
    
```

Reducer

< Rec, null-value >

---

---

---

---

---

---

---

---

CentraleSupélec

## Patron de partitionnement – v2

< Offset, Record >

```

PartAtt = Record.extract("Year")
if (PartAtt < Y1)
  Partid = 0
else if (PartAtt < Y2)
  Partid = 1
...
write < Partid, Record >
    
```

Mapper

Hyp: Partition map: [Y<sub>1</sub>, Y<sub>2</sub>, ..., Y<sub>n-1</sub>] (n parts)

< Partid, Record >

```

Partitioner(key (= Partid)) {
  return(key)
}
    
```

Partitioner

Hyp: Nb of Reducers = Nb of parts (n)

Shuffle & Sort

< Partid, [Rec1, Rec2...] (= ListVal) >

```

Hyp: Nb of Reducers = Nb of Parts
for each record in ListVal
  write < null-key, record >
  Or:
  write < record, null-value >
    
```

Reducer

< Rec, null-value >

---

---

---

---

---

---

---

---

**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*

---

---

---

---

---

---

---

---

---

---

**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<sup>ème</sup> 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

---

---

---

---

---

---

---

---

---

---

**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*

**1<sup>ère</sup> étape : partitionnement**  
 $E_0 : x < s_1$   
 $E_1 : s_1 \leq x < s_2$   
 $E_2 : s_2 \leq x < s_3$   
 $E_3 : s_3 \leq x$

**2<sup>ème</sup> étape : tri de chaque sous-ensemble**  
 $E_0 : x < s_1$   
 $E_1 : s_1 \leq x < s_2$   
 $E_2 : s_2 \leq x < s_3$   
 $E_3 : s_3 \leq x$

File-0  
File-1  
File-2  
File-3

- Extension du patron de partitionnement : ajoute un tri local dans chaque *Reducer*
- Besoin de connaissance a priori sur les données pour définir la carte de partitionnement équilibrée (les seuils)

---

---

---

---

---

---

---

---

CentraleSupélec

## Tri total avec tri final dans les *Reducers*

```

< Offset, Record >
  Val = Record.extract("Att1")
  if (Val < S1)
    Partid = 0
  else if (Val < S2)
    Partid = 1
  ...
  write < Partid, Record >
  
```

Partition map (thresholds) : {S<sub>1</sub>, S<sub>2</sub>, ..., S<sub>n-1</sub>}

Mapper  
Using the Partition map  
(n parts, n-1 thresholds)

```

< Partid, Record >
  Partitioner(key (= Partid)) {
    return(key)
  }
  
```

Partitioner  
Hyp: Nb of Reducers = Nb of parts (= n)

Shuffle & Sort

```

< Partid, [Rec1, Rec2...] (= ListVal) >
  // Sort records, function of their "Att1" attribute
  ListVal = sort(ListVal)
  // Write the sorted list of records
  for each record in ListVal
    write < record, null-value >
  
```

Reducer  
Hyp: Nb of Reducers = Nb of Parts (= n)

< Rec, null-value >

---

---

---

---

---

---

---

---

### Tri total avec tri final dans les *Reducers*

**Déploiement et bilan**

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

```

    graph LR
      subgraph Mappers
        M1[Mapper]
        M2[Mapper]
        M3[Mapper]
      end
      subgraph Partitioners
        P1[Partitioner]
        P2[Partitioner]
        P3[Partitioner]
      end
      subgraph ShuffleSort
        SS[Shuffle & Sort]
      end
      subgraph Reducers
        R1[Reducer]
        R2[Reducer]
      end
      M1 --> P1
      M2 --> P2
      M3 --> P3
      P1 --> SS
      P2 --> SS
      P3 --> SS
      SS --> R1
      SS --> R2
  
```

- 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*...**

---

---

---

---

---

---

---

---

### 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

---

---

---

---

---

---

---

---

### 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 exæquos selon k2 en ordre décroissant, puis ...

---

---

---

---

---

---

---

---

### Reconfiguration du Shuffle & Sort

1<sup>er</sup> niveau de contrôle/reconfiguration : **redéfinition du Partitioner**

Un mapper → <k, v> → ? → [Reducer 0, Reducer 1, Reducer 2, ..., Reducer n-1]

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

---

---

---

---

---

---

---

---

### Reconfiguration du Shuffle & Sort

1<sup>er</sup> 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);
    
```

---

---

---

---

---

---

---

---

### Reconfiguration du Shuffle & Sort

2<sup>ème</sup> niveau de reconfiguration : **redéfinition du keyComparator**  
 (ou sortComparator)

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

Un mapper → <k3, v3> → keyComparator → Reducer i  
 Un mapper → <k2, v2> → keyComparator → Reducer i  
 Un mapper → <k4, v4> → keyComparator → Reducer i

keyComparator → Reducer i

[ <k1, v1>  
 <k2, v2>  
 <k3, v3>  
 <k4, v4>  
 <k5, v5> ]

---

---

---

---

---

---

---

---

**Reconfiguration du Shuffle & Sort**

2<sup>ème</sup> niveau de reconfiguration : **redéfinition du *keyComparator***  
(ou *sortComparator*)

Hyp : clés composites  
 $k1 < k2 < k3$   
 $t11 < t12$ , et  $t21 < t22$

Un mapper < {k2,t21}, v3 >  
 Un mapper < {k1,t12}, v2 > < {k3,t31}, v5 >  
 Un mapper < {k2,t22}, v4 >

keyComparator sur les 2 attributs de la clé composite

---

---

---

---

---

---

---

---

---

---

**Reconfiguration du Shuffle & Sort**

2<sup>ème</sup> 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);
    
```

---

---

---

---

---

---

---

---

---

---

**Reconfiguration du Shuffle & Sort**

2<sup>ème</sup> niveau de reconfiguration : **redéfinition du *groupComparator***

Hyp : clés composites  
 $k1 < k2 < k3$   
 $t11 < t12$ , et  $t21 < t22$

Un mapper < {k2,t21}, v3 >  
 Un mapper < {k1,t12}, v2 > < {k3,t31}, v5 >  
 Un mapper < {k2,t22}, v4 >

groupComparator sur 1<sup>er</sup> attribut de la clé seulement

keyComparator  
groupComparator

Reducer i

reduce(< {k1,\_}, [v1, v2] >)  
 reduce(< {k2,\_}, [v3, v4] >)  
 reduce(< {k3,\_}, [v5] >)

---

---

---

---

---

---

---

---

---

---

## Reconfiguration du *Shuffle & Sort*

2<sup>ème</sup> niveau de reconfiguration : **redéfinition du *groupComparator***

```

public static class MyGroupComparator
  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 groupComparator to the Map-Reduce job
job.setGroupingComparatorClass(MyGroupComparator.class)
// Run the Map-Reduce job
System.exit(job.waitForCompletion(true) ? 0 : 1);
    
```

---

---

---

---

---

---

---

---

---

---

## 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

---

---

---

---

---

---

---

---

---

---

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

**Mapper**

```

for each < Offset, Record >
  V = Record.extract("Att1")
  write < V, Record >
    
```

**Partitioner:**  
 Hyp: nb of Reducers = nb of part. blocks

**Partition map (thresholds):**  
 $\{s_1, s_2, \dots, s_{n-1}\}$

- if (key <  $s_1$ ) return 0
- else if (key <  $s_2$ ) return 1
- .....
- else return (n-1)

**Shuffle & Sort**

**keyComparator:**  
 (controls the ASC or DESC order)

- if ( $v_1 < v_2$ ) return (-/+1)
- else if ( $v_1 > v_2$ ) return (+/-1)
- else return(0)

**Reducer**  
 Hyp: nb of Reducers = nb of parts (n)

```

for each Rec in ListVal
  write < Rec, null-value >
    
```

---

---

---

---

---

---

---

---

---

---



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

Déploiement et bilan

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

```

    graph LR
      M1[Mapper] --> P1[Partitioner]
      M2[Mapper] --> P2[Partitioner]
      M3[Mapper] --> P3[Partitioner]
      P1 --> SS[Shuffle & Sort]
      P2 --> SS
      P3 --> SS
      SS --> KC[keyComparator]
      KC --> R1[Reducer]
      KC --> R2[Reducer]
      KC --> R3[Reducer]
  
```

- 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

---

---

---

---

---

---

---

---

**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

---

---

---

---

---

---

---

---

**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*

---

---

---

---

---

---

---

---

## 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

---

---

---

---

---

---

---

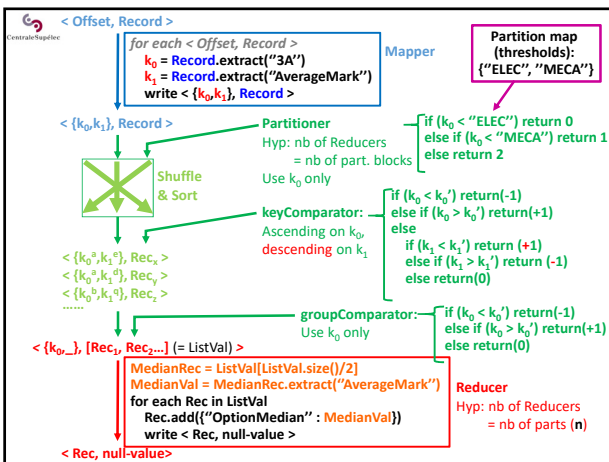
---

---

---

---

---




---

---

---

---

---

---

---

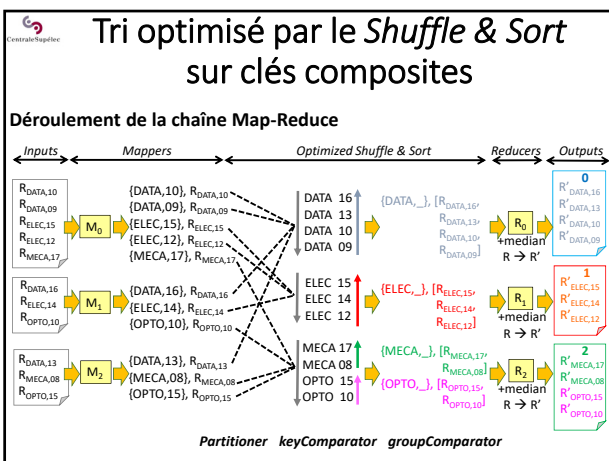
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

---

---

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

**Déploiement et bilan**

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

Mapper → Partitioner → Shuffle & Sort → keyComparator → groupComparator → Reducer

- Tri multi-critères (emboîtés)
- Exploite les tris de clés du *Shuffle & Sort*, **aucun tri final dans les Reducers**
- Regroupement dans les *Reducers* et traitements par sous-listes
- **Suppose l'existence d'un plan de partitionnement prédéfini...**

---

---

---

---

---

---

---

---

---

---

---

---

### Patrons de jointure

1 – Jointure de 2 tables dans 2 fichiers distribués

2 – Jointure de grande taille avec prétraitement

---

---

---

---

---

---

---

---

---

---

---

---

### Jointure de 2 tables dans 2 fichiers distribués

JOIN A.Att3 = B.AttZ

Result

```
{Att1 : 10,
Att2 : 100,
Att3 : 128,
.....
AttZ : 128,
AttY : 'bob',
AttX : .....
}

{Att1 : 10,
Att2 : 100,
Att3 : 128,
.....
AttZ : 128,
AttY : 'ben',
AttX : ....
}
```

---

---

---

---

---

---

---

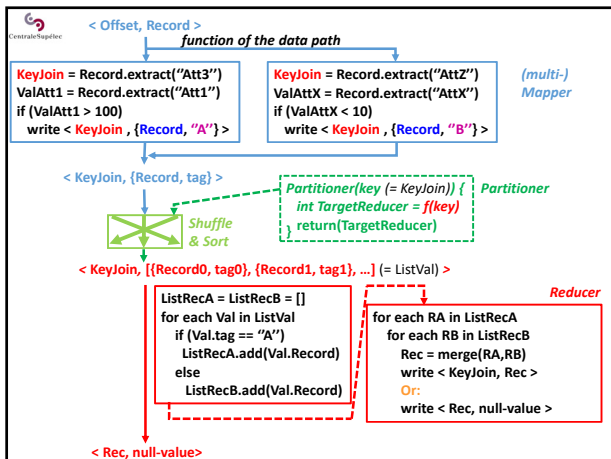
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

### Jointure de 2 tables dans 2 fichiers distribués

**1 - Paire {clé, liste de valeurs} reçue par un Reducer :**

```

< KeyJoin = "Dupont", [{R1,"A"}, {R2,"A"}, {R'1,"B"}, {R'2,"B"}, {R3,"A"}] >
    
```

**2 - Actions du Reducer :**

**Re-construction de deux listes d'enregistrements :**

```

→ ListRecA = (R1, R2, R3)
   ListRecB = (R'1, R'2)
    
```

**Génération de paires clé – valeur de jointure**

```

→ < "Dupont", R1 U R'1 >   ou bien : → < R1 U R'1, null-value >
   < "Dupont", R1 U R'2 >   < R1 U R'2, null-value >
   < "Dupont", R2 U R'1 >   < R2 U R'1, null-value >
   < "Dupont", R2 U R'2 >   < R2 U R'2, null-value >
   < "Dupont", R3 U R'1 >   < R3 U R'1, null-value >
   < "Dupont", R3 U R'2 >   < R3 U R'2, null-value >
    
```

---

---

---

---

---

---

---

---

---

---

### Jointure de 2 tables dans 2 fichiers distribués

**Déploiement:**

Prévoir beaucoup de Reducers, car le résultat du Join peut être gros, et entraîner :

- beaucoup de traitements dans les Reducers
- beaucoup d'E/S sur disques!

---

---

---

---

---


---

---

---

---

---

 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

---

---

---


---

---

---

---

---

 CentraleSupélec

## 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

---

---

---

---

---


---

---

---

---

---

 **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

---

---

---

---

---


---

---

---

---

---

 **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);
    
```

---

---

---

---

---

---

---

---

---

---

**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);
    
```

Definir : un job Map-Reduce  
un partitionner  
un échantillonneur

↓

Les configurer et les associer

Exécuter l'échantillonneur qui génère  
la carte de partitionnement

↓

Exécuter le job Map-Reduce avec le  
partitionneur qui exploite la carte

---

---

---

---

---

---

---

---

---

---

**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 générique
- 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



---

---

---

---

---

---

---

---