




Big Data
Technologies Internes d'Hadoop

Stéphane Vialle & Gianluca Quercini




1



Principes et technologie d'Hadoop

1. **Localité des données et des traitements**
2. Framework d'Hadoop
3. Mécanismes du Map-Reduce d'Hadoop
4. Système de fichiers distribué d'Hadoop (HDFS)
5. Allocation et gestion de ressources d'Hadoop

2



Localité des données et des traitements

C'est toujours l'accès aux données qui coute cher, pas le calcul lui-même (une fois les données arrivées dans l'unité de calcul)

Big Data *façon Hadoop* :
 Amener les codes de traitements aux données
 → Transformer momentanément en nœuds de traitement les nœuds de stockage des données traitées
 → Eviter de déplacer des données (très volumineuses) ... mais ...
 → les relire et les réécrire localement chaque fois que la RAM est pleine

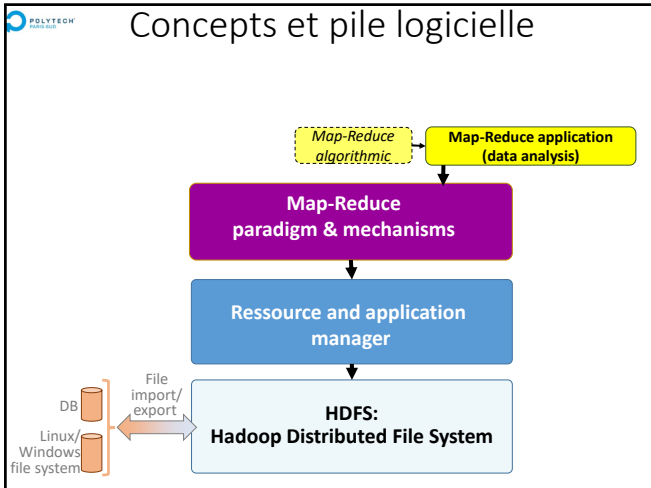
Big Data *façon Spark* :
 Amener les données à toutes les unités de calcul disponibles
 → Mais ... lire les données et les écrire sur disque une seule fois
 → et ... garder les données en RAM durant tout le traitement

3

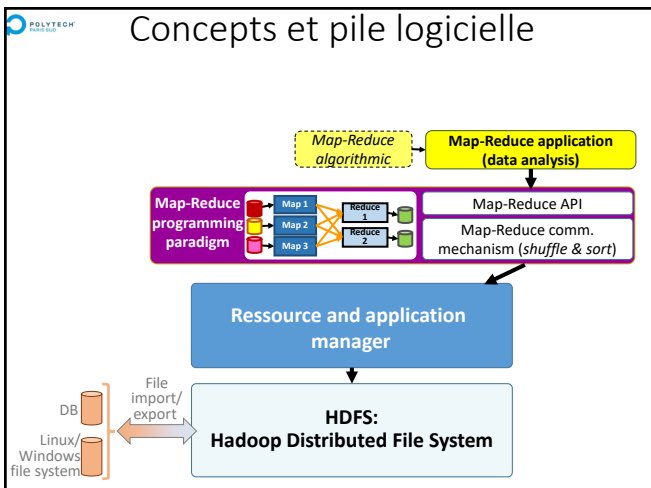
Principes et technologie d'Hadoop

1. Localité des données et des traitements
2. **Framework d'Hadoop**
3. Mécanismes du Map-Reduce d'Hadoop
4. Système de fichiers distribué d'Hadoop (HDFS)
5. Allocation et gestion de ressources d'Hadoop

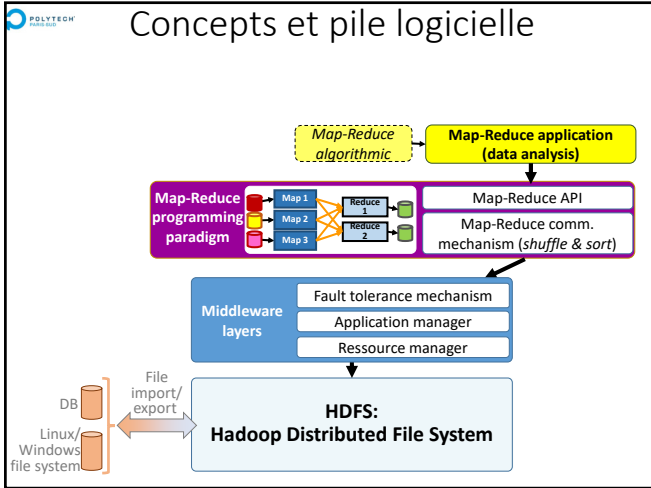
4



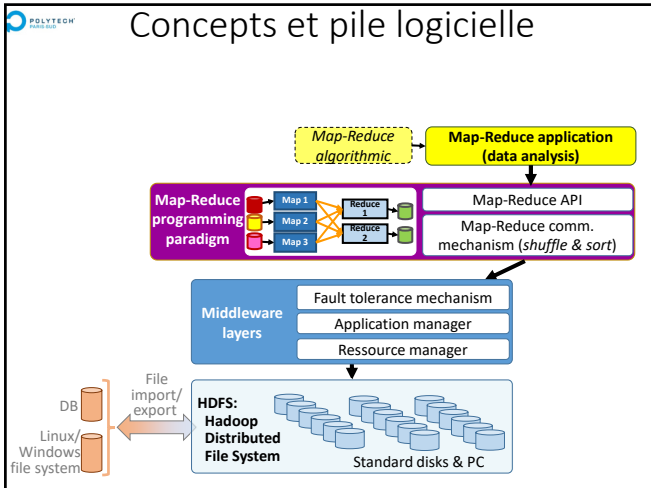
5



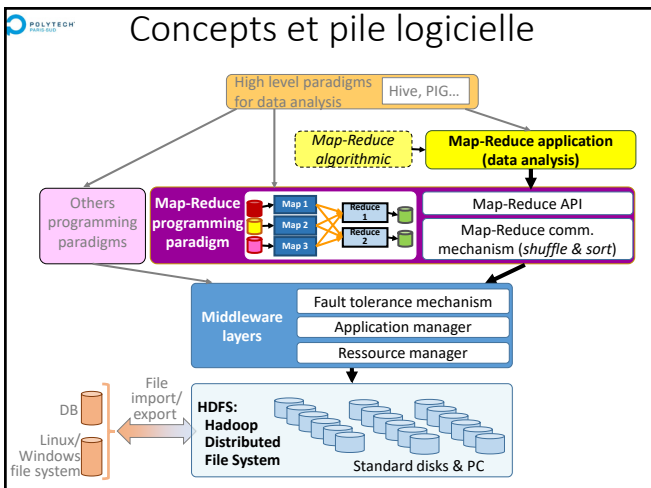
6



7



8



9

Principes et technologie d'Hadoop

1. Localité des données et des traitements
2. Framework d'Hadoop
- 3. Mécanismes du Map-Reduce d'Hadoop**
4. Système de fichiers distribué d'Hadoop (HDFS)
5. Allocation et gestion de ressources d'Hadoop

10

Chaîne d'opérations Map-Reduce

(17001, Pignon)
(16995, Martin)
(17012, Durand)
(16200, Dupont)
(17003, Dupond)
(16158, Martin)

Map function

```
f_map(key,val) {
  if key < 17000
    write (val, 1)
}
```

(Martin, 1)
(Dupont, 1)
(Martin, 1)
(Martin, (1,1))
(Dupont, (1))

Reduce function

```
g_reduce(key, list_vals) {
  sum = 0
  for each val in list_vals
    sum += val
  write (key, sum)
}
```

Suite de transformations de paires « clé-valeur »

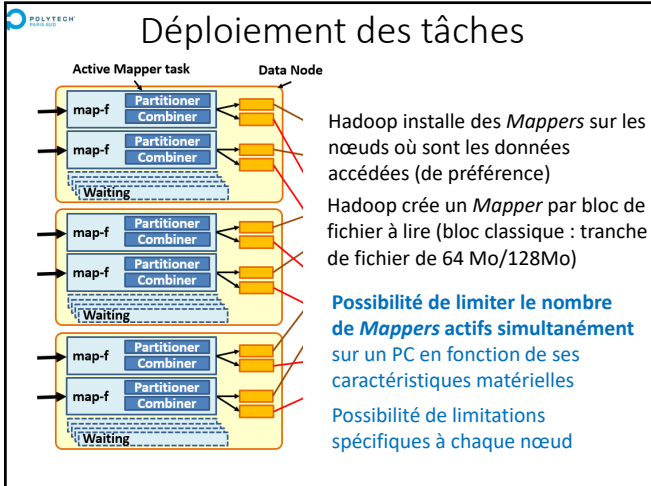
11

Chaîne d'opérations Map-Reduce

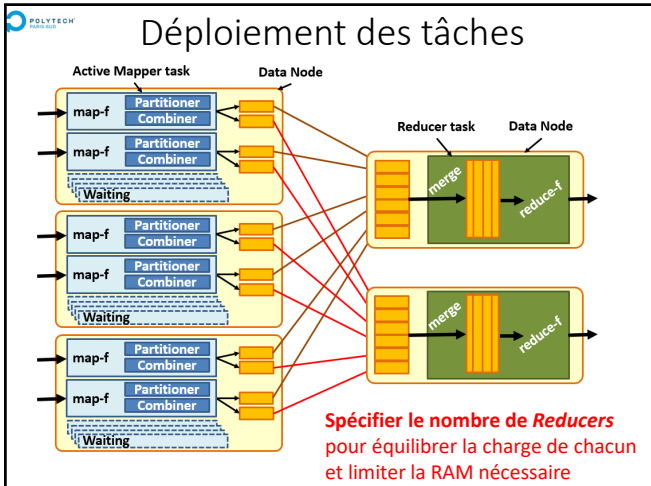
Input Data → **Record Reader** (Input data reading from HDFS (customizable)) → **Records (key/values)₁** → **Mapper** (Treatment 1: data filtering & first local computations) → **Records (key/values)₂** → **Combiner** (Treatment 2: in order to reduce communications (optional & usually close to Reducer code)) → **Records (key/values)₃** → **Partitioner** (Associate keys to Reducers. Default version can be optimized for the application.) → **Shuffle and Sort** (Key-Value pair routage. Predefined comm. and sorting scheme) → **Records (key/[val])₃** → **Reducer** (Treatment 3: local computations on data grouped with identical keys) → **Records (key/val)₄** → **Output Format** (Output result writing in HDFS (customizable)) → **Output Data**

keyComparator(k1,k2)
groupComparator(k1,k2)
Key comparison functions, impacting order and grouping of key-value pairs supplied to Reducers. Default functions can be overloaded to adapt to the application.

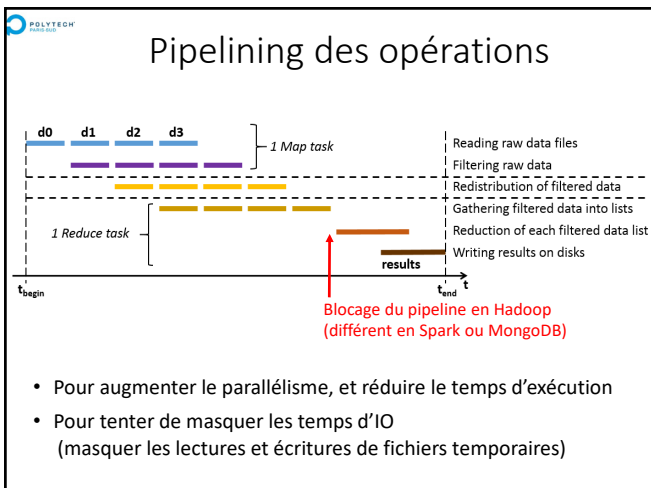
12



13



14



15

Principes et technologie d'Hadoop

1. Localité des données et des traitements
2. Framework d'Hadoop
3. Mécanismes du Map-Reduce d'Hadoop
4. **Système de fichiers distribué d'Hadoop (HDFS)**
 - Distribution et réplication des fichiers sous HDFS
 - Lecture de fichiers sous HDFS
 - Ecriture de fichiers sous HDFS
5. Allocation et gestion de ressources d'Hadoop

16

Mécanismes d'HDFS

Le *NameNode* conserve la cartographie du HDFS + les évolutions (les « logs »)
 → Permet de savoir où sont les fichiers

Le *secondary NameNode* stocke les logs et met à jour la cartographie (qd bcp de logs)
 → Sur un nœud à part pour pouvoir calculer la mise à jour sans ralentir le système

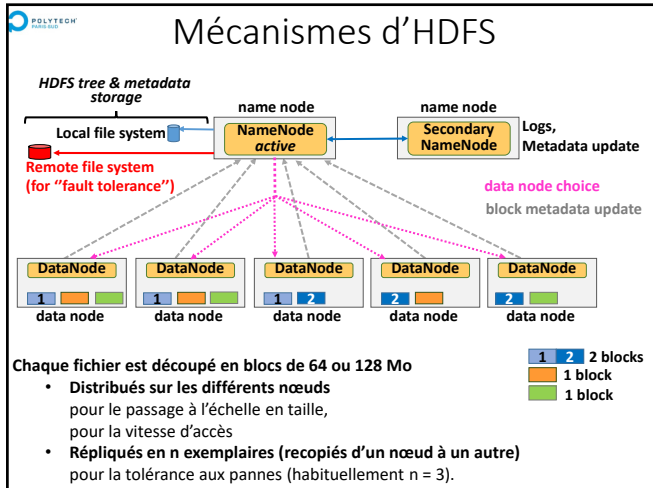
17

Mécanismes d'HDFS

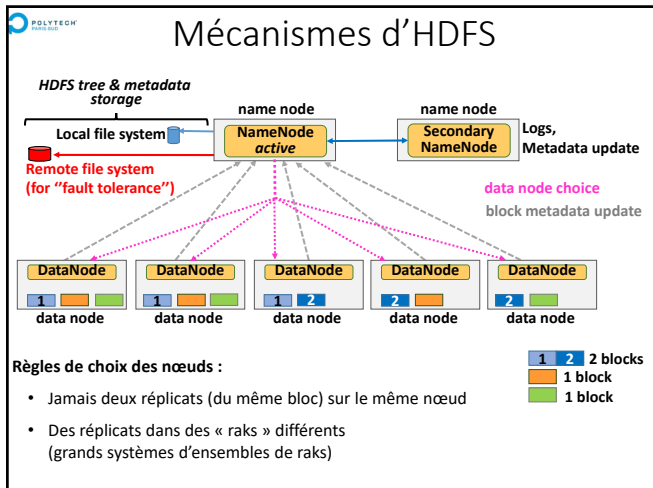
Les métadonnées du HDFS sont stockés sur le file system classique, localement.

Un stockage distant permet de renforcer la tolérance aux pannes
 (sans ses métadonnées, le HDFS est inexploitable)

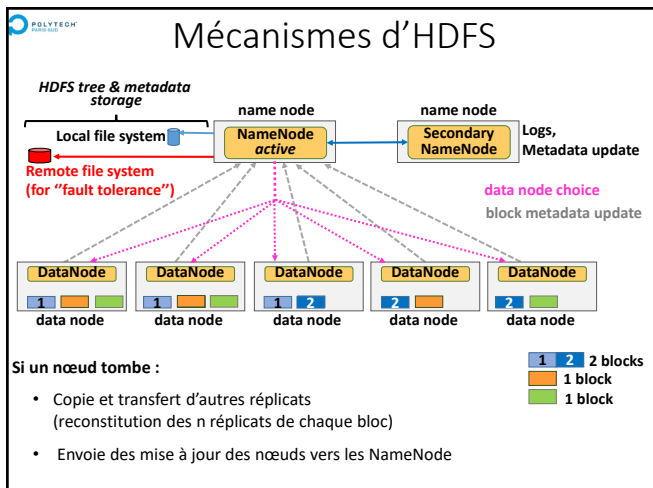
18



19



20



21

Mécanismes d'HDFS

Pourquoi des blocs de 64 Mo ?

Temps de « seek » : temps de positionnement au début du fichier sur le disque (disque standard, rotatif)
 $T_{seek} = 10ms$

Bande passante disque std : $Bw = 100 Mo/s$
 $Tread = Q / Bw$

On veut : $T_{seek} < 1\% Tread$
 $\Leftrightarrow 10 \cdot 10^{-3} s < (1/100) \cdot (Q/100) s$
 $\Leftrightarrow 100 (Mo) < Q$

→ Des blocs de 64 Mo ou 128 Mo permettent de masquer les temps de seek

→ Au-delà : pas plus de gain, mais moins de distribution des fichiers, moins de vitesse de lecture

22

Mécanismes d'HDFS : « haute disponibilité »

Le NameNode est un SPOF : Single Point Of Failure
 → Si on perd le NameNode alors le File System d'Hadoop ne marche plus !

23

Mécanismes d'HDFS : « haute disponibilité »

Si on perd le NameNode alors le File System d'Hadoop ne marche plus !
 → On duplique le NameNode

24

Mécanismes d'HDFS : « haute disponibilité »

The diagram shows two NameNodes: one labeled 'NameNode active' and another 'Secondary NameNode'. They are connected by a bidirectional arrow. Below them are five DataNodes. The first three DataNodes have a '1' in a blue box, and the last two have a '2' in a blue box. Dashed lines connect the active NameNode to all DataNodes, and the Secondary NameNode to the DataNodes with '2'.

- Le NameNode en *standby* est passé *actif* « très rapidement »
- Il n'y a plus de SPOF

→ On « ne sent plus passer la panne » : Haute Disponibilité

25

Principes et technologie d'Hadoop

1. Localité des données et des traitements
2. Framework d'Hadoop
3. Mécanismes du Map-Reduce d'Hadoop
4. **Système de fichiers distribué d'Hadoop (HDFS)**
 - Distribution et réplication des fichiers sous HDFS
 - **Lecture de fichiers sous HDFS**
 - Ecriture de fichiers sous HDFS
5. Allocation et gestion de ressources d'Hadoop

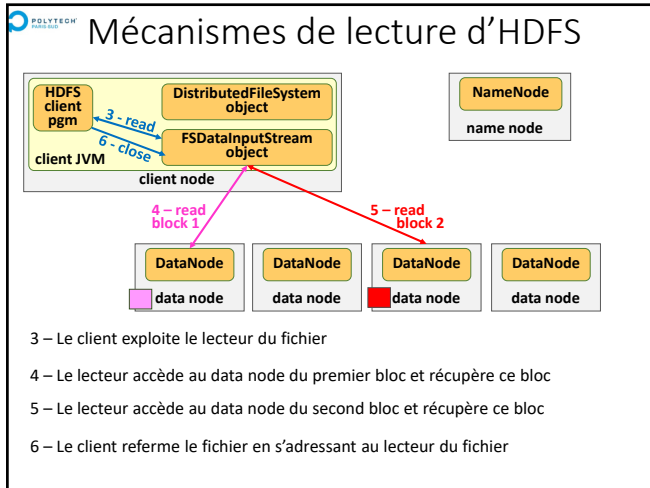
26

Mécanismes de lecture d'HDFS

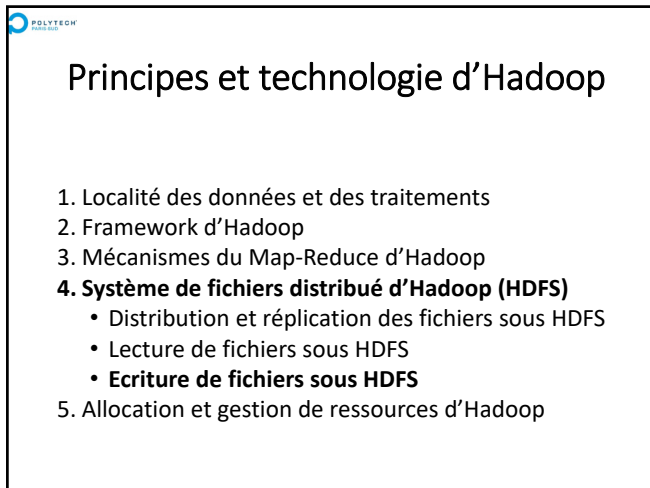
The diagram shows a 'client node' containing 'HDFS client pgm', 'client JVM', 'DistributedFileSystem object', and 'FSDataInputStream object'. It interacts with a 'NameNode' (name node). Arrows indicate: '1 - open' from client to DistributedFileSystem; '2a - get block locations' from DistributedFileSystem to NameNode; and '2b - create object' from NameNode to FSDataInputStream.

- 0 – Création d'un objet permettant de s'interfacer à HDFS (un stub/proxy d'HDFS)
- 1 – Demande d'ouverture d'un fichier HDFS en lecture (« open »)
- 2 – Demande de localisation du fichier
 - 2a - Le proxy interroge le NameNode : pour savoir quels blocs lire et où les lire
Le NameNode répond au proxy
 - 2b - Le proxy crée et retourne un objet lecteur spécialisé sur le fichier ciblé

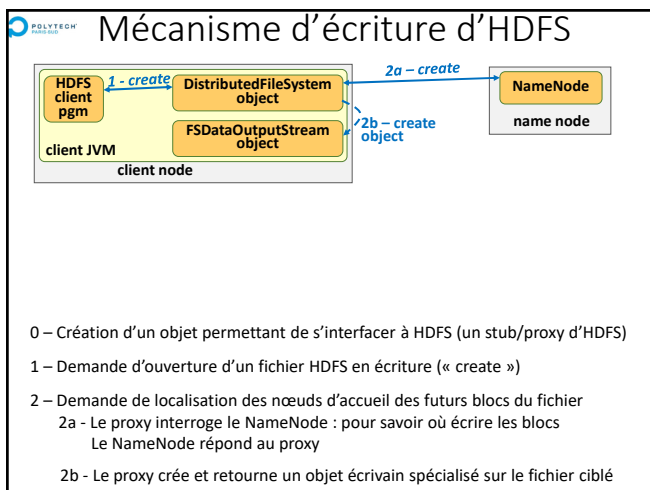
27



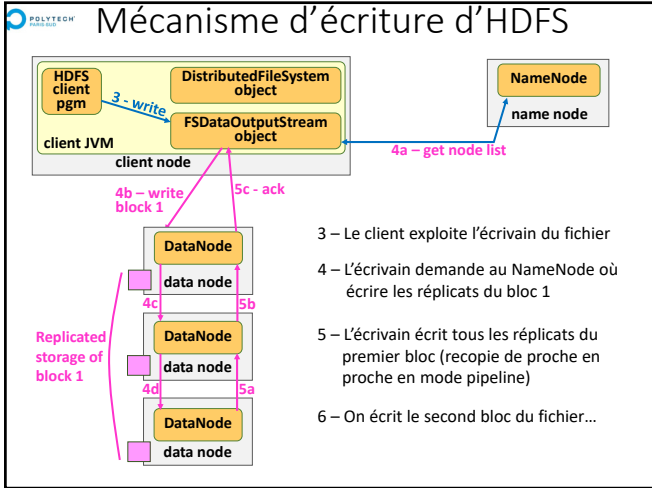
28



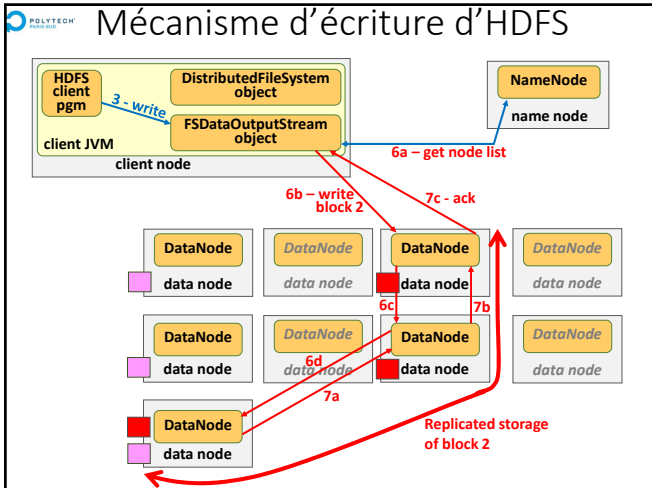
29



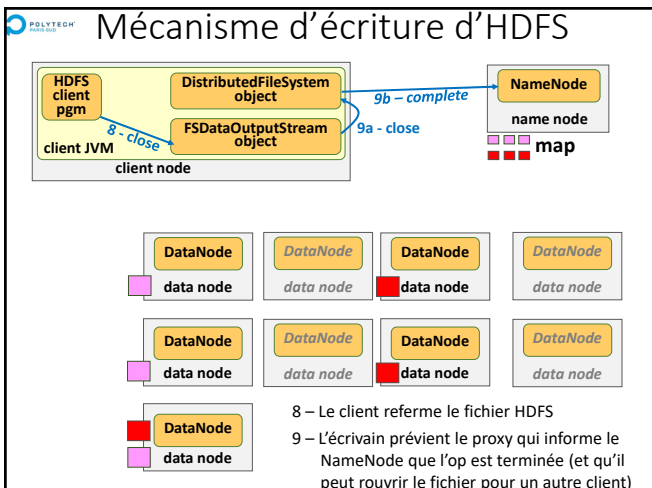
30



31



32



33

Principes et technologie d'Hadoop

1. Localité des données et des traitements
2. Framework d'Hadoop
3. Mécanismes du Map-Reduce d'Hadoop
4. Système de fichiers distribué d'Hadoop (HDFS)
- 5. Allocation et gestion de ressources d'Hadoop**
 - Hadoop v1
 - Hadoop v2 (YARN) : passage à l'échelle amélioré

34

Gestion des ressources v1

The client creates a proxy local object to ensure all communications with Hadoop rsrc & job manager, and with HDFS services

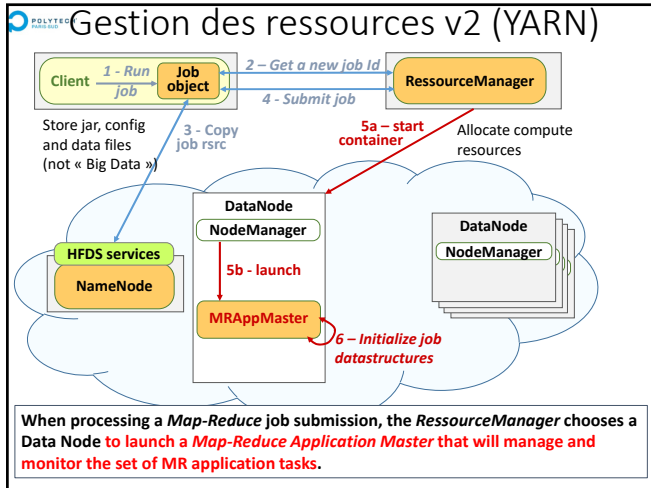
35

Gestion des ressources v1

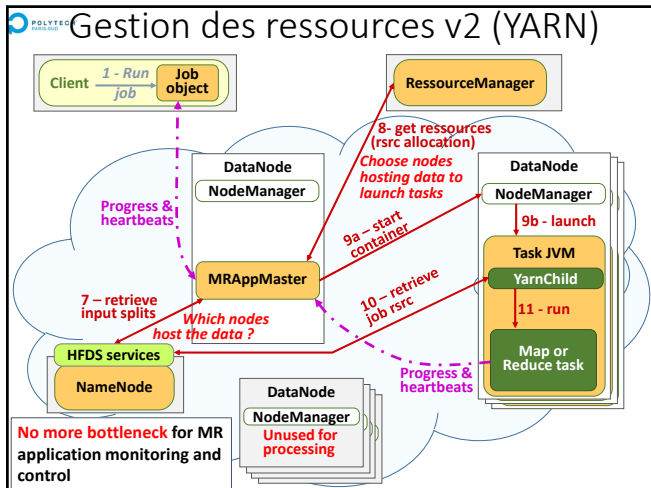
1 - Run job
2 - Get a new job id
3 - Copy job rsrc (not « Big Data »)
4 - Submit job
5 - Initialize job datastructures
6 - Retrieve input splits: find the data nodes to contact

The JobTracker asks to HDFS NameNode where are stored the target data: in order to choose the data nodes that will support Map and Reduce tasks →

36



40



41

Exécution spéculative

Hadoop lance de nombreuses tâches (map, reduce, ...). Certaines peuvent « planter » ou « ralentir ».

Le TaskTracker (MR v1) ou l'ApplicationManager (MR v2) monitore fortement les exécutions des tâches, et détectent ces « ralentissements ».


Hadoop peut faire une exécution spéculative : il lance de nouvelles instances des tâches « en retard », et dès qu'il obtient des résultats d'une tâche, il tue son doublon.

Mais cette démarche a un coût :

- en charge de calcul (création fréquentes de tâches redondantes)
- en déplacement de données (surtout si on redonne un reducer)

On peut débrayer ce comportement (ex : cluster HPC très fiable)


42

 QUIZ

Q1: a Hadoop Map-Reduce program ends up generating a huge number of key-value pairs with (always) the same key

- Will the HDFS output file be stored in one large block or in several small ones?
- Will the “reduce” treatment be processed in parallel or sequentially?


43

 QUIZ

Q2: a user connects his client program, running on his laptop, to a 100-node Hadoop cluster, and submits Map-Reduce queries, to compute the histogram of the age of the French (with one-year increments)

- Technically, can he download the results to his laptop?
- Technically, can he upload new input data to the HDFS of the 100-node cluster?
- Technically, can he download the input data to his laptop and then load it into a second Hadoop cluster?
- Is it possible to copy data from the HDFS of a first Hadoop cluster directly to the HDFS of a second?

44

 QUIZ

Q3: a failure occurs on a Hadoop data node used during the execution of a Map-Reduce program (the node disappears)

- Does the user have to resubmit the Map-Reduce request?
- Does the user get the result later when a failure occurs?


Q4: to improve fault tolerance, you can install HDFS on top of a RAID-enabled storage array (*Redundant Array of Independent Disks*)

- Do you think this is a logical approach?

45

POLYTECH

Technologies Internes d'Hadoop



46
