

UNIVERSITÉ PARIS-SACLAY POLYTECH PARIS-SACLAY

Big Data

Technologies Internes d'Hadoop

Stéphane Vialle & Gianluca Quercini

université PARIS-SACLAY
CentraleSupélec

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

L1SN
LABORATOIRE NATIONAL D'INFORMATIQUE DES SCIENCES DU NUMÉRIQUE

Grand Est
Région Île de France

1

POLYTECH PARIS-SACLAY

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

POLYTECH PARIS-SACLAY

Localité des données et des traitements

C'est toujours l'accès aux données qui coûte 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

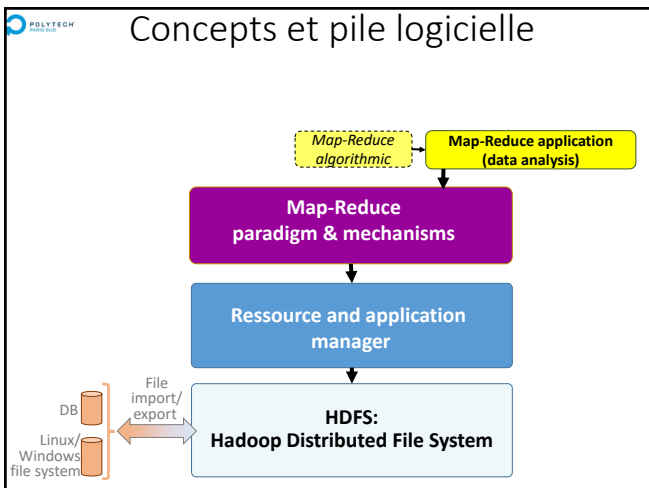
3

POLYTECH PARIS-SACLAY

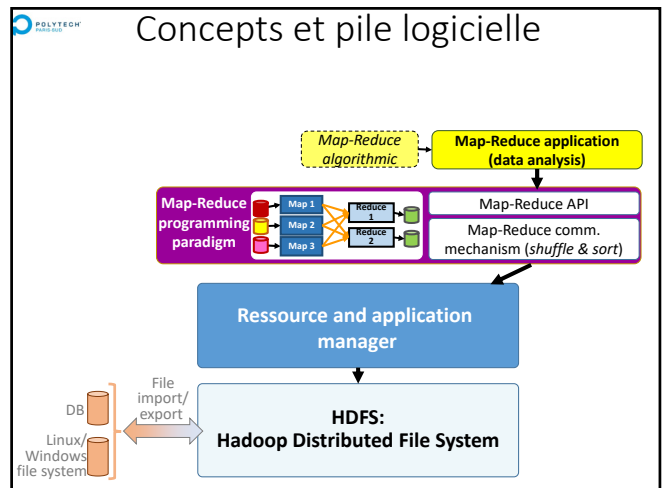
Principes et technologie d'Hadoop

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

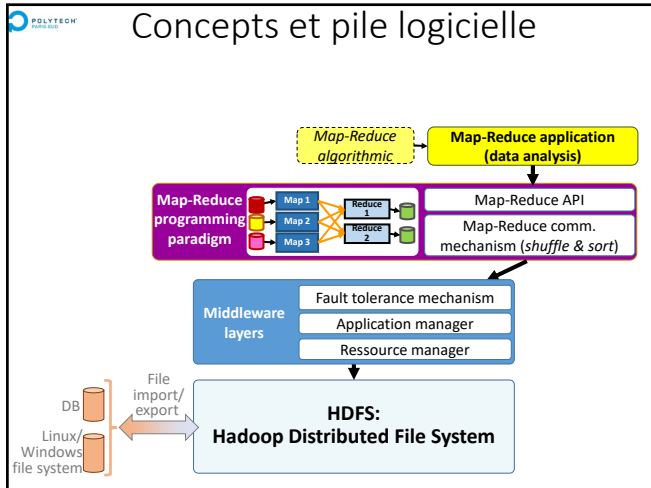
4



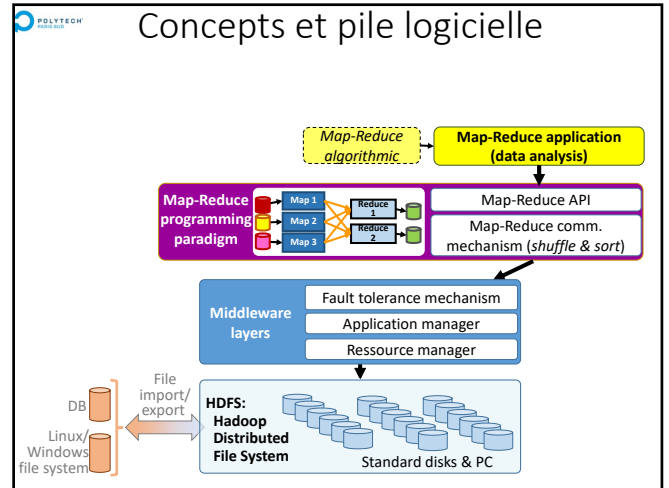
5



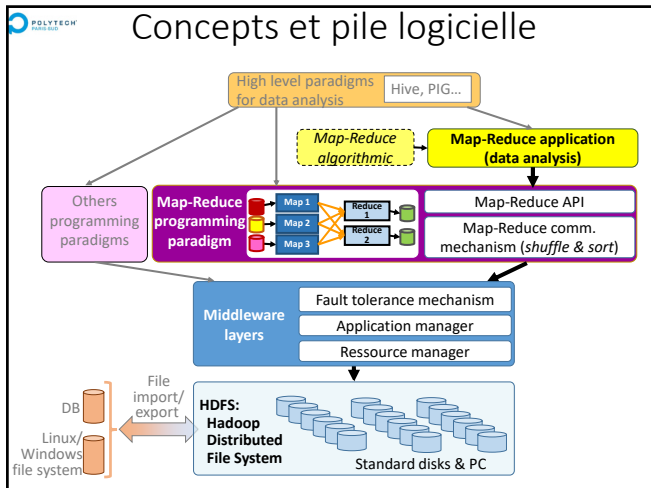
6



7



8

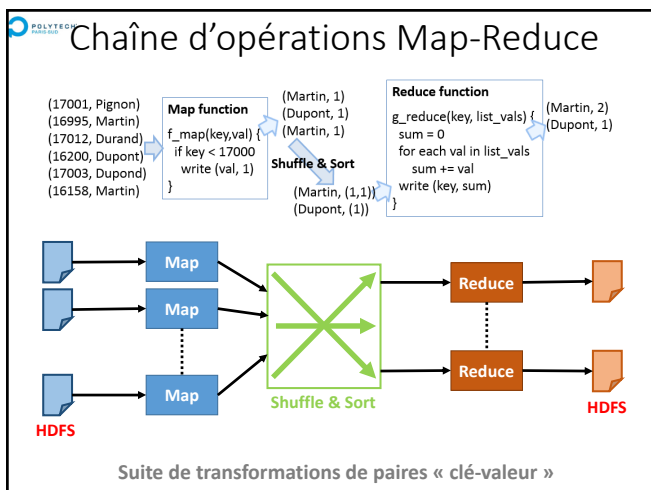


9

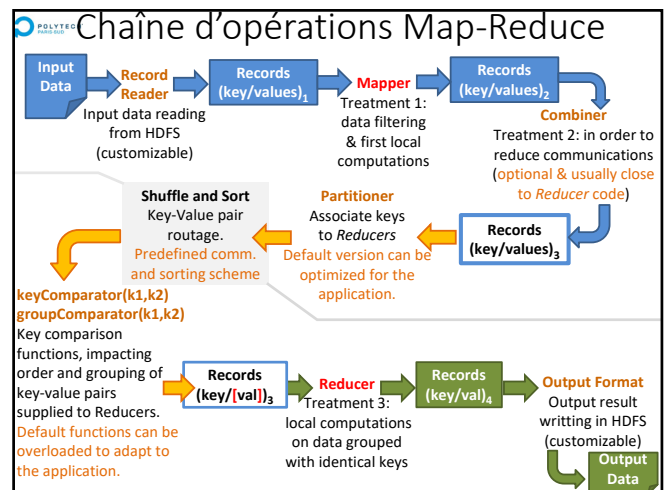
Principes et technologie d'Hadoop

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

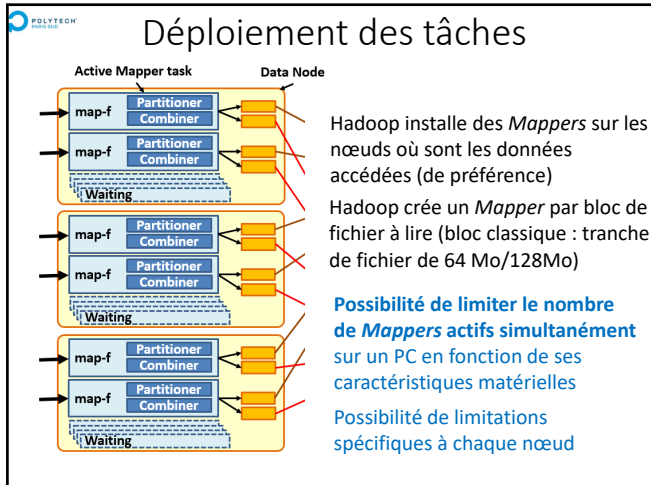
10



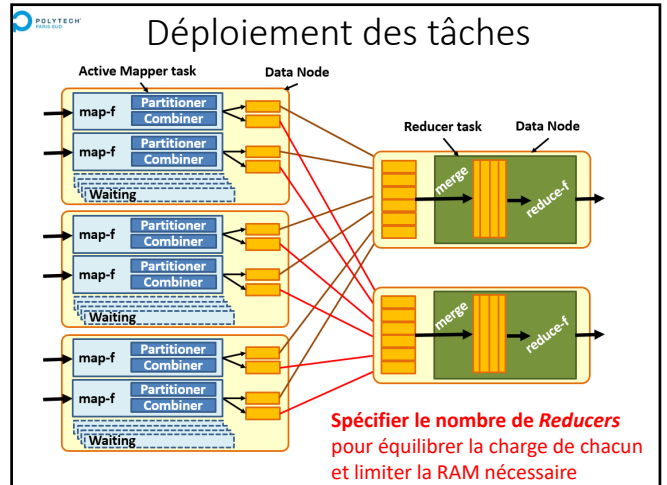
11



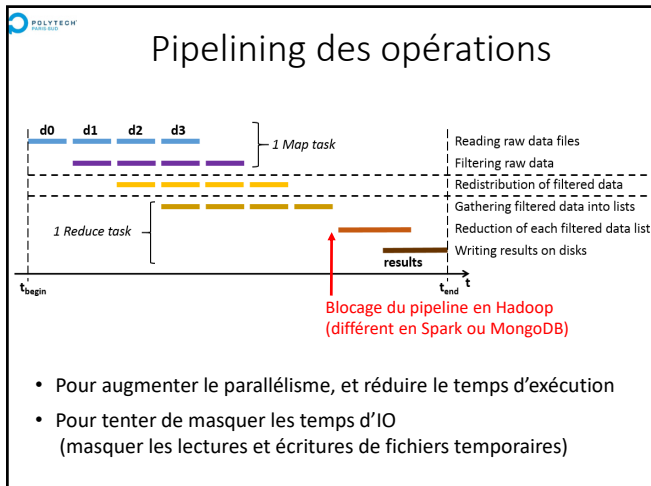
12



13



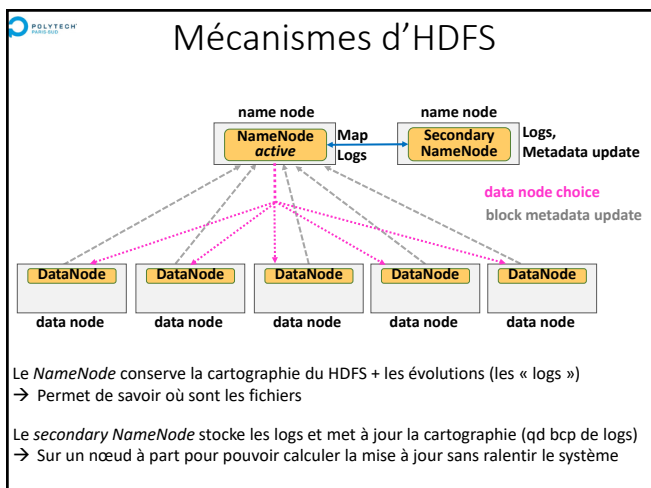
14



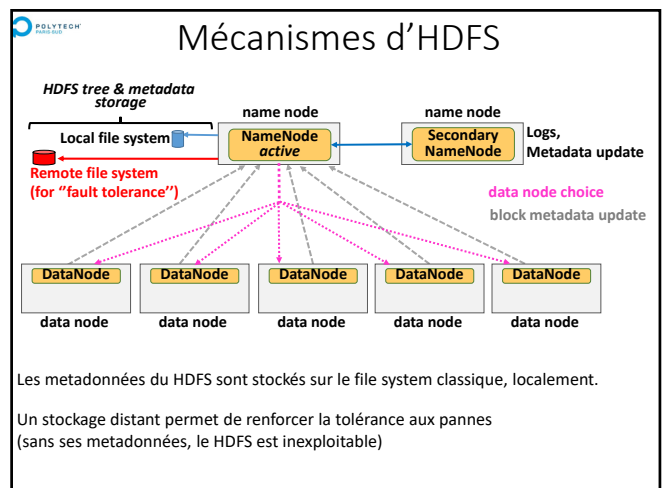
15

- ### 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



17



18

Mécanismes d'HDFS

HDFS tree & metadata storage

Local file system → NameNode active

Remote file system (for "fault tolerance")

Logs, Metadata update

data node choice
block metadata update

Chaque fichier est découpé en blocs de 64 ou 128 Mo

- Distribués sur les différents nœuds pour le passage à l'échelle en taille, pour la vitesse d'accès
- Répliqués en n exemplaires (recopiés d'un nœud à un autre) pour la tolérance aux pannes (habituellement $n = 3$).

1 2 2 blocks
1 block
1 block

19

Mécanismes d'HDFS

HDFS tree & metadata storage

Local file system → NameNode active

Remote file system (for "fault tolerance")

Logs, Metadata update

data node choice
block metadata update

Règles de choix des nœuds :

- Jamais deux répliquats (du même bloc) sur le même nœud
- Des répliquats dans des « racks » différents (grands systèmes d'ensembles de racks)

1 2 2 blocks
1 block
1 block

20

Mécanismes d'HDFS

HDFS tree & metadata storage

Local file system → NameNode active

Remote file system (for "fault tolerance")

Logs, Metadata update

data node choice
block metadata update

Si un nœud tombe :

- Copie et transfert d'autres répliquats (reconstitution des n répliquats de chaque bloc)
- Envoie des mise à jour des nœuds vers les NameNode

1 2 2 blocks
1 block
1 block

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 »

SPOF

name node
NameNode active

name node
Secondary NameNode

name node
NameNode standby

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 »

name node
NameNode active

name node
Secondary NameNode

name node
NameNode standby

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

→ On duplique le NameNode

24

Mécanismes d'HDFS : « haute disponibilité »

- 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

- 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

Mécanismes de lecture d'HDFS

- 3 – Le client exploite le lecteur du fichier
- 4 – Le lecteur accède au data node du premier bloc et récupère ce bloc
- 5 – Le lecteur accède au data node du second bloc et récupère ce bloc
- 6 – Le client referme le fichier en s'adressant au lecteur du fichier

28

Principes et technologie d'Hadoop

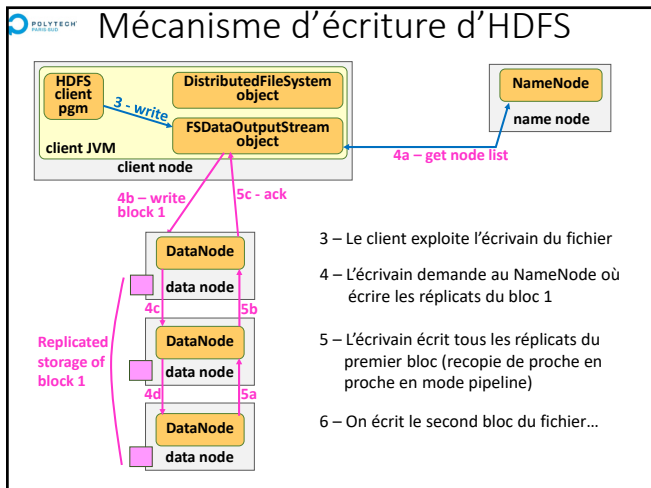
1. Localité des données et des traitements
2. Framework d'Hadoop
3. Mécanismes du Map-Reduce d'Hadoop
- 4. Système de fichiers distribué d'Hadoop (HDFS)**
 - Distribution et réplication des fichiers sous HDFS
 - Lecture de fichiers sous HDFS
 - **Ecriture de fichiers sous HDFS**
5. Allocation et gestion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

29

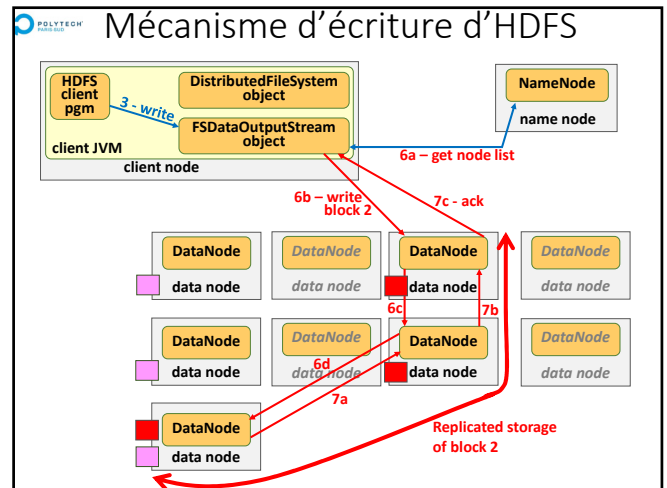
Mécanisme d'écriture d'HDFS

- 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 écriture (« create »)
- 2 – Demande de localisation des nœuds d'accueil des futurs blocs du fichier
 - 2a - Le proxy interroge le NameNode : pour savoir où écrire les blocs
Le NameNode répond au proxy
 - 2b - Le proxy crée et retourne un objet écrivain spécialisé sur le fichier ciblé

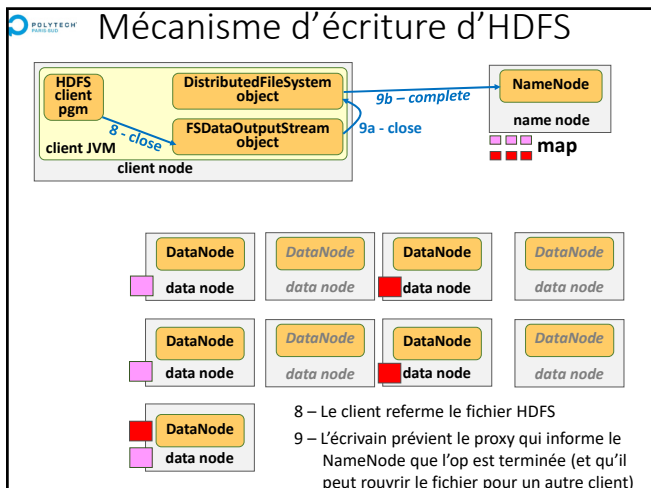
30



31



32

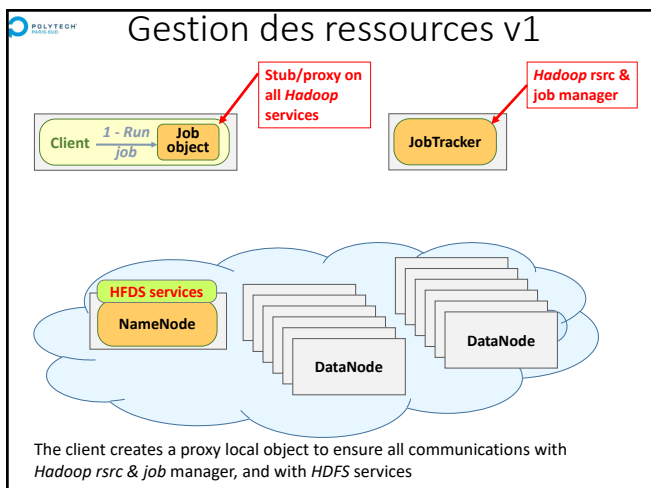


33

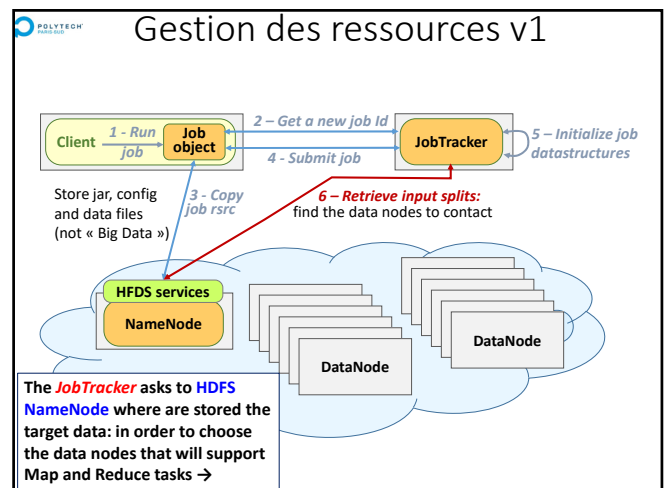
Principes et technologie d'Hadoop

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

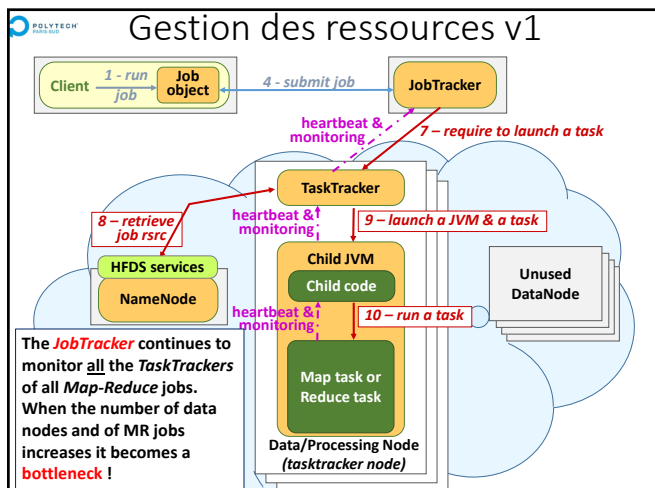
34



35



36

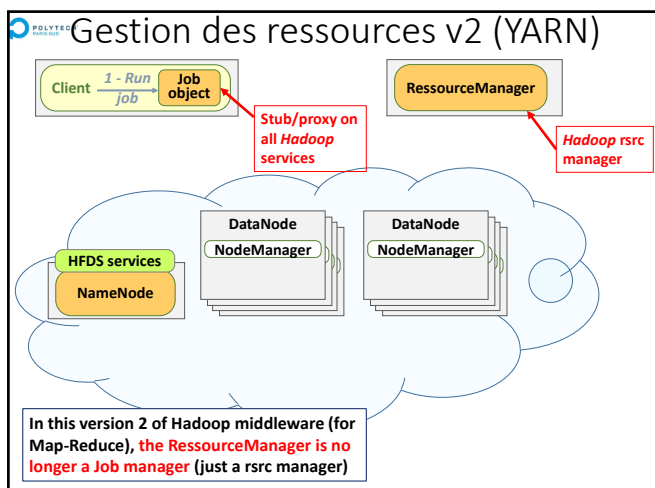


37

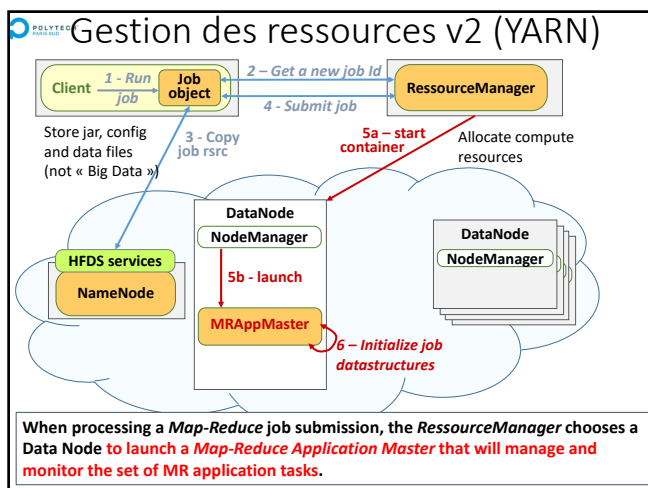
Principes et technologie d'Hadoop

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**

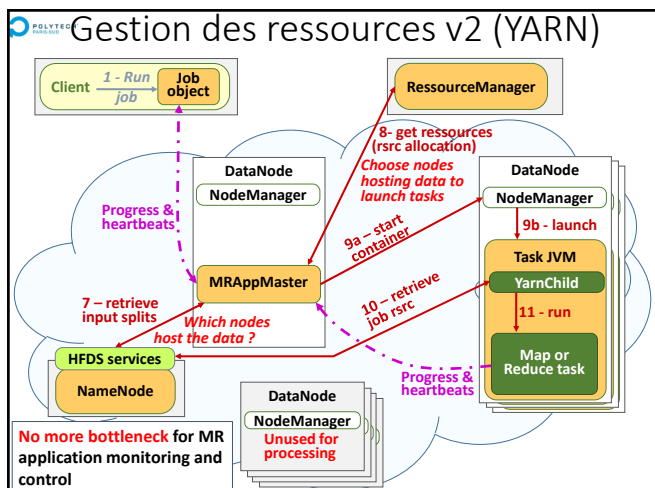
38



39



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) monitoré 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

POLYTECH
UNIVERSITÄT
DUISBURG
ESSEN

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

POLYTECH
UNIVERSITÄT
DUISBURG
ESSEN

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

POLYTECH
UNIVERSITÄT
DUISBURG
ESSEN

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
UNIVERSITÄT
DUISBURG
ESSEN

Technologies Internes d'Hadoop



46