

Big Data

Principes d'HDFS

Stéphane Vialle

&

Gianluca Quercini



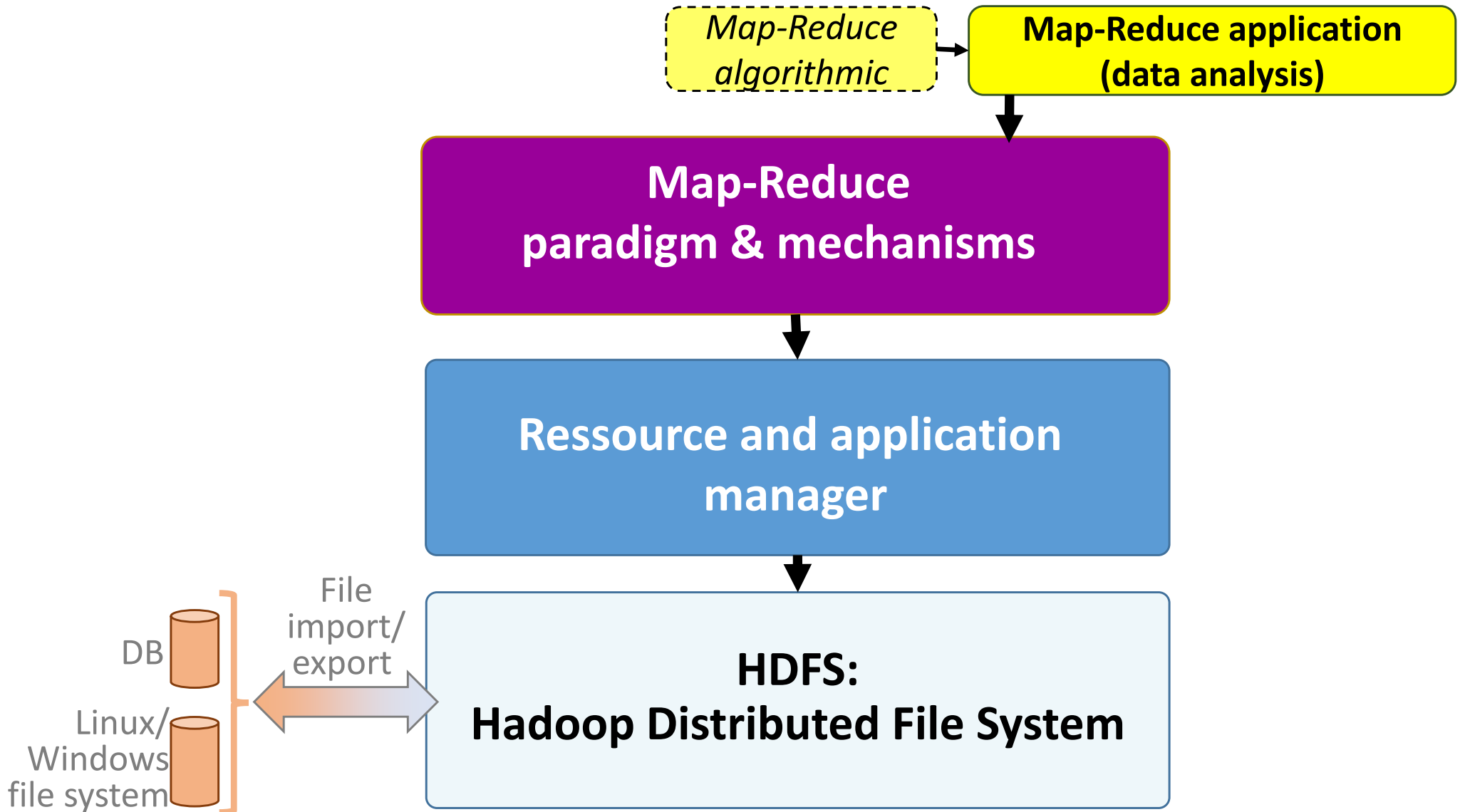
ÉCOLE DOCTORALE
Sciences et technologies
de l'information
et de la communication (STIC)



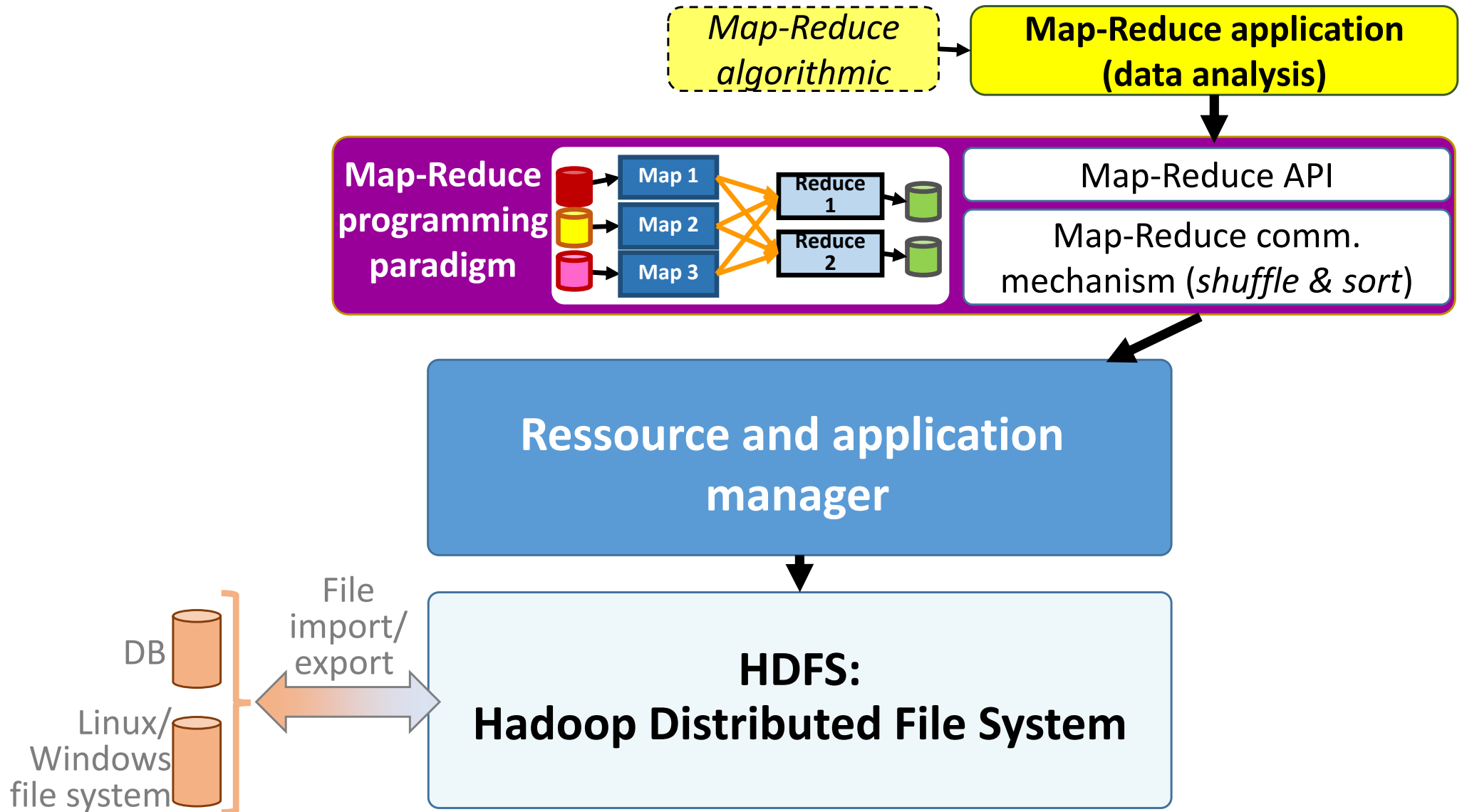
Principes d'HDFS

- 1. Framework d'Hadoop**
2. Système de fichiers distribué d'Hadoop : HDFS

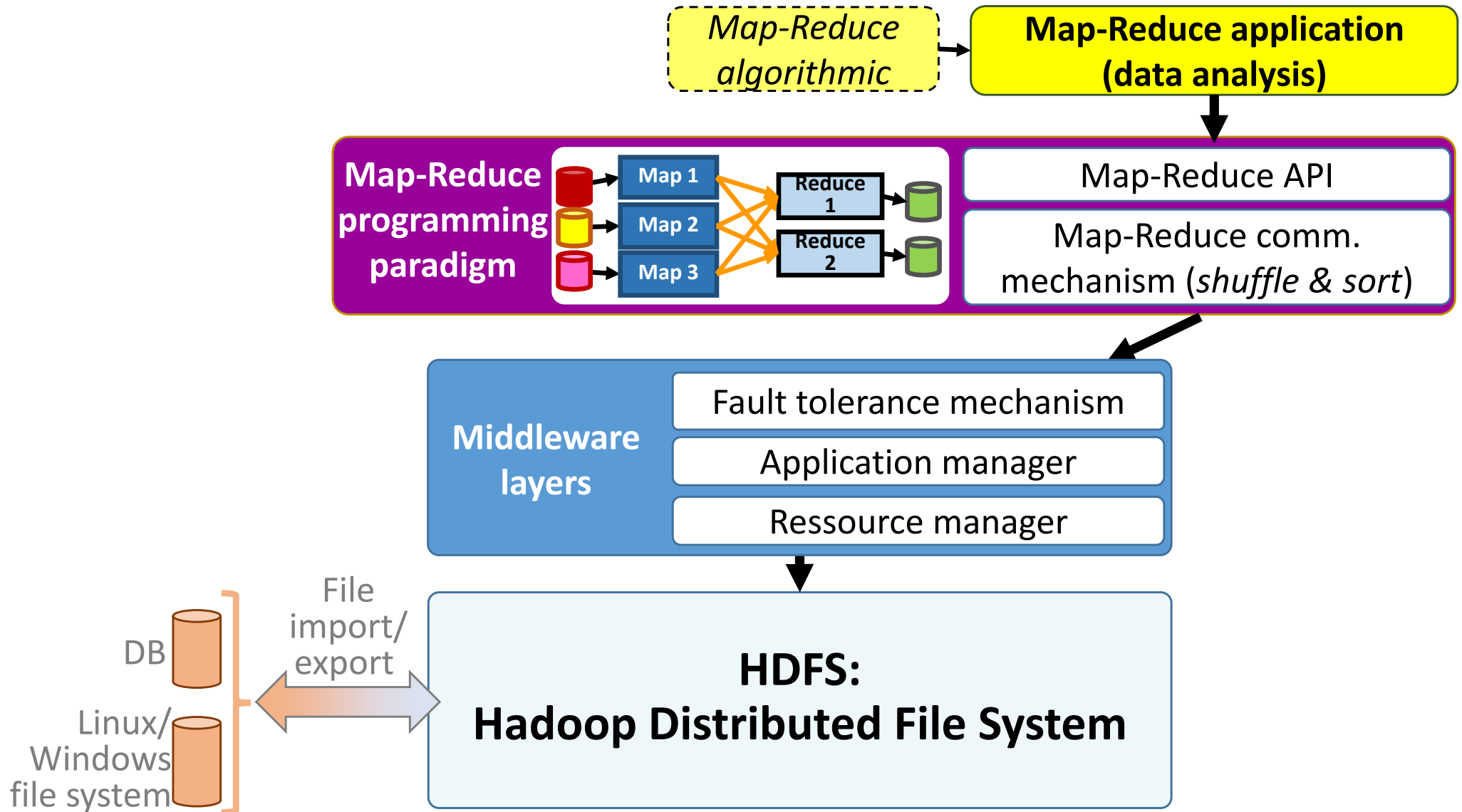
Concepts et pile logicielle



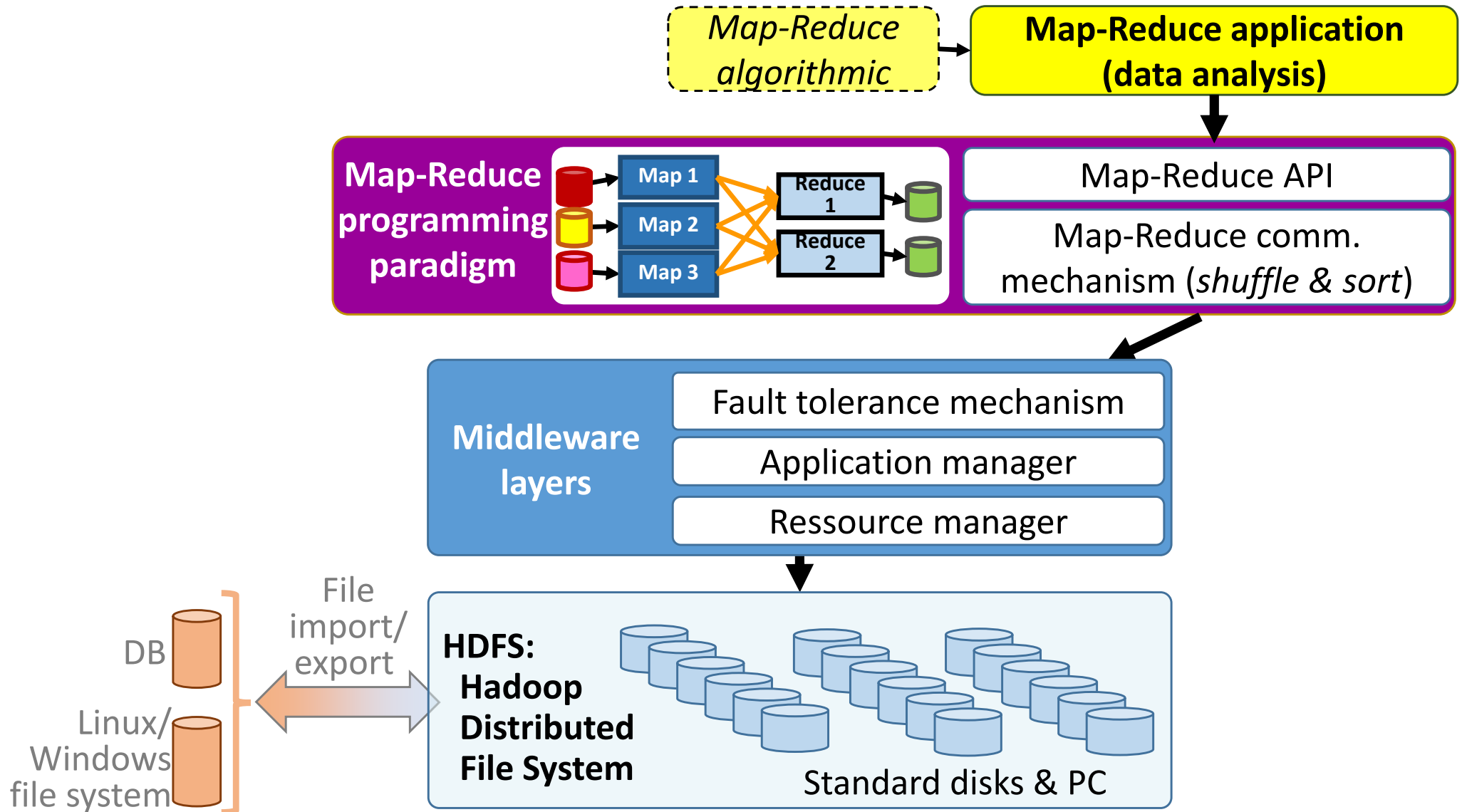
Concepts et pile logicielle



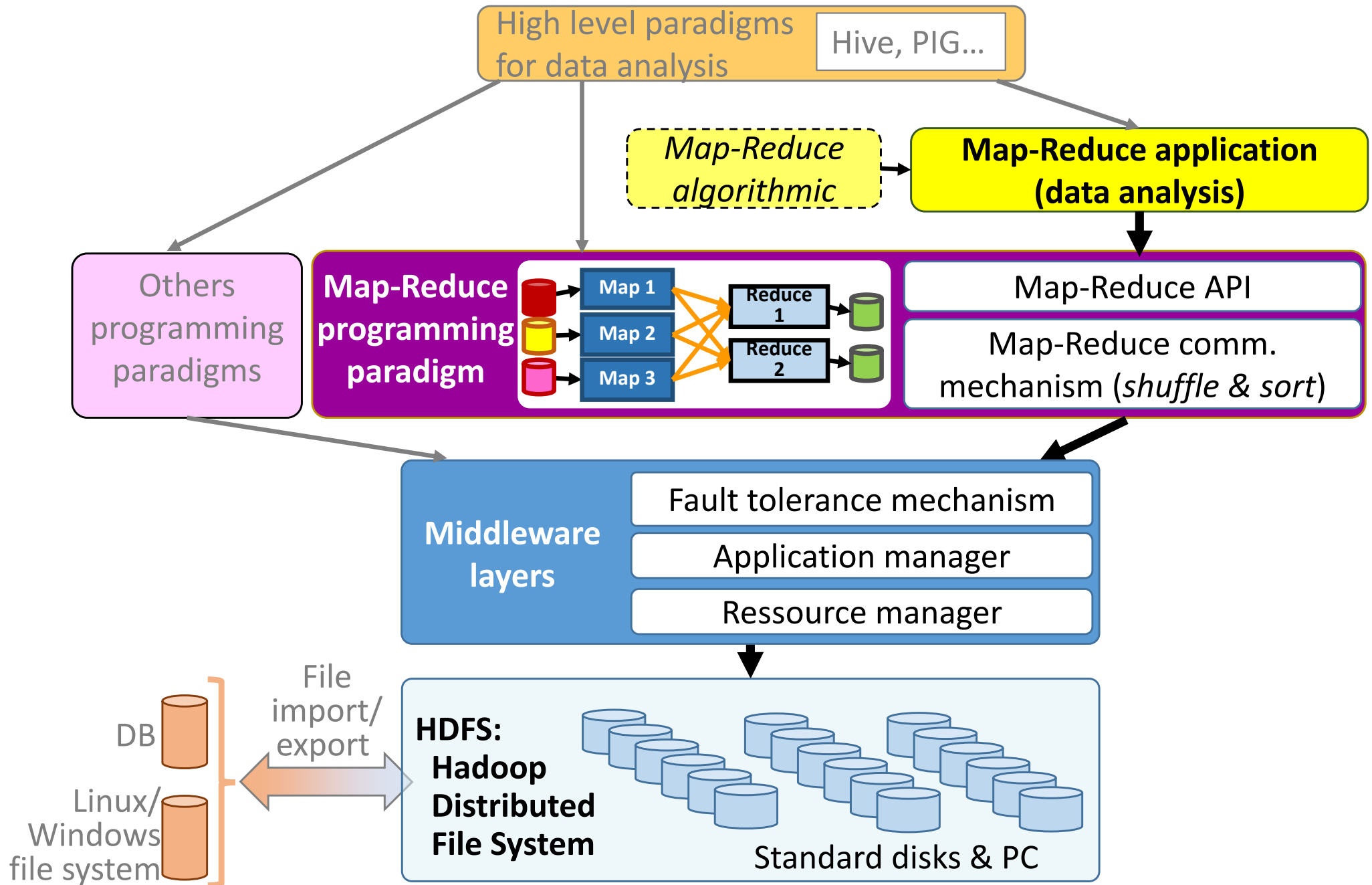
Concepts et pile logicielle



Concepts et pile logicielle



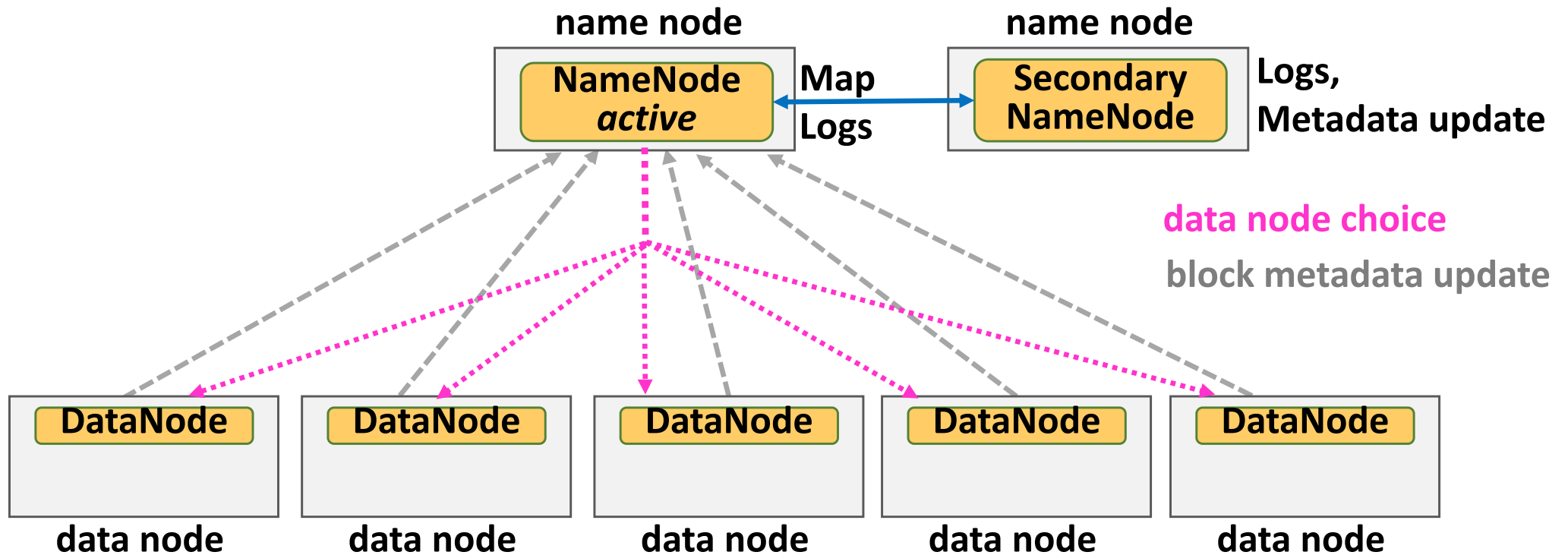
Concepts et pile logicielle



Principes d'HDFS

1. Framework d'Hadoop
- 2. HDFS (Hadoop Distributed File System)**
 - 2.1. Stockage des fichiers sous HDFS**
 - 2.2. Mécanisme de haute disponibilité d'HDFS
 - 2.3. Lecture de fichiers sous HDFS
 - 2.4. Ecriture de fichiers sous HDFS

Mécanismes d'HDFS



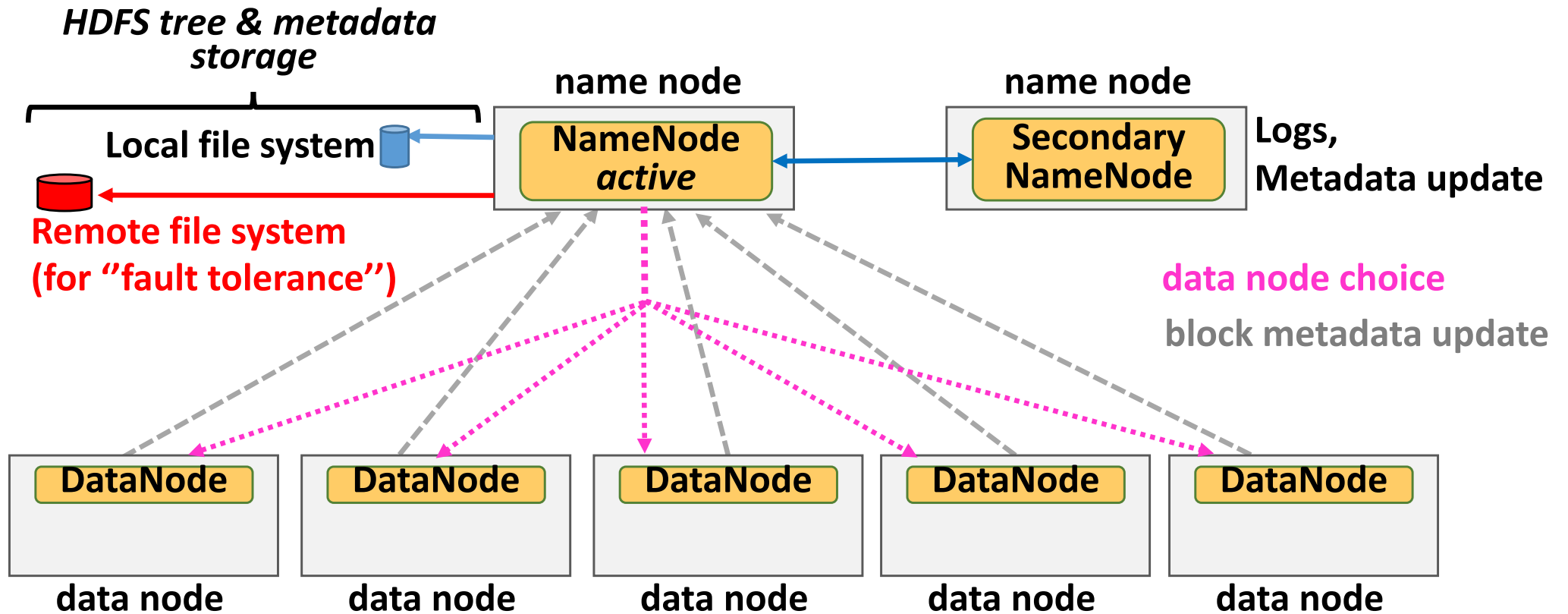
Le *NameNode* conserve la cartographie du HDFS + les évolutions (les « logs »)

→ 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

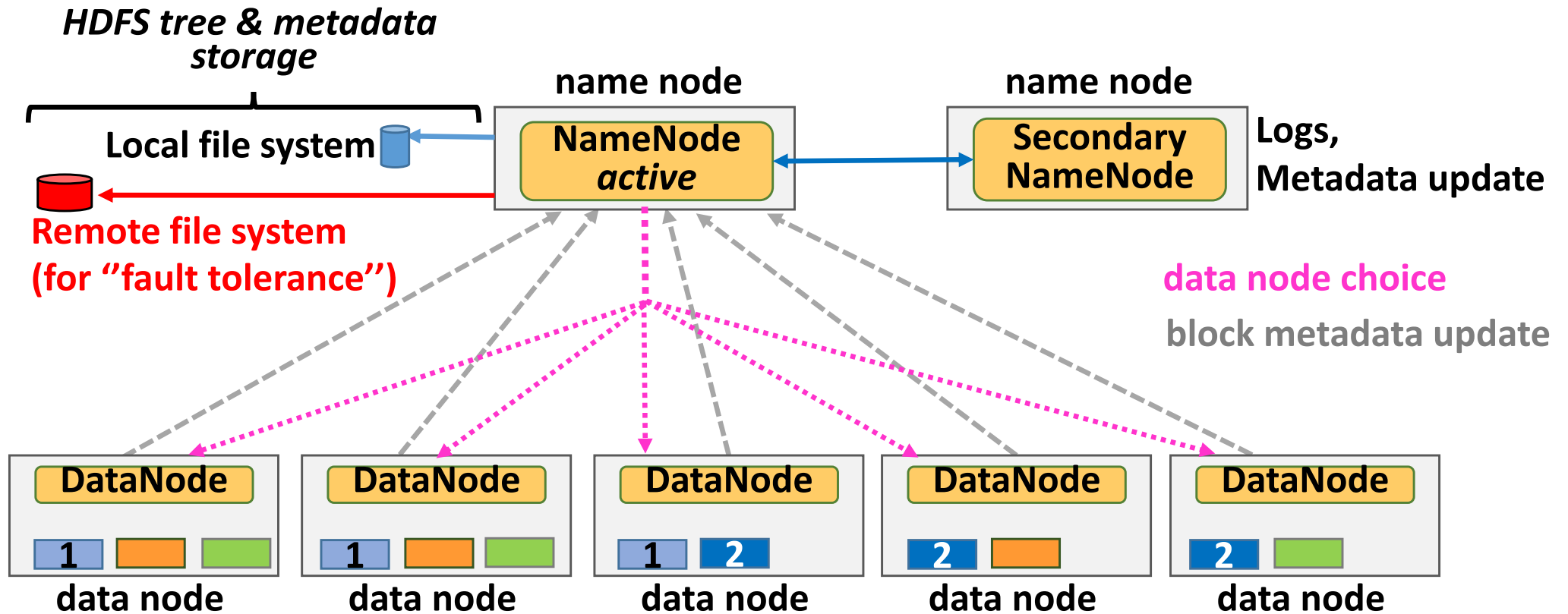
Mécanismes d'HDFS



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

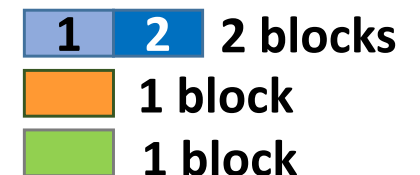
Un stockage distant permet de renforcer la tolérance aux pannes (sans ses métadonnées, le HDFS est inexploitable)

Mécanismes d'HDFS

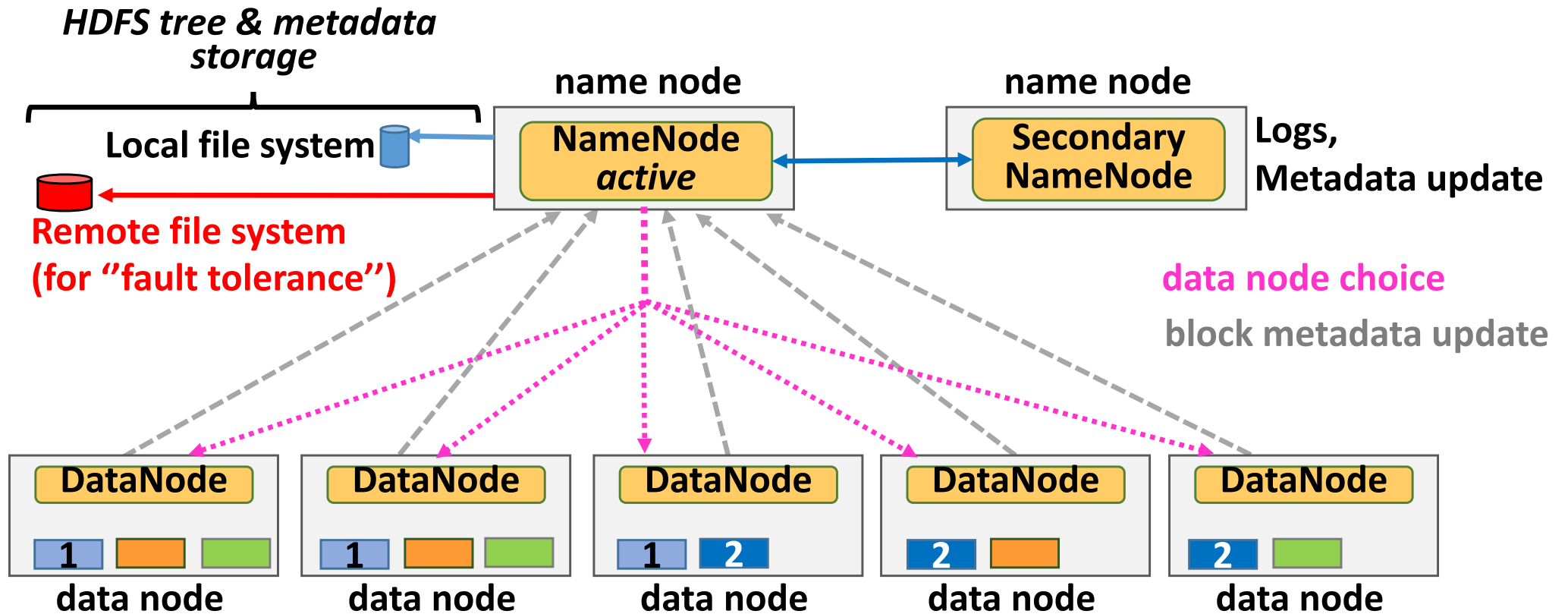


Chaque fichier est découpé en blocs de 64 ou 128 Mo

- 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).

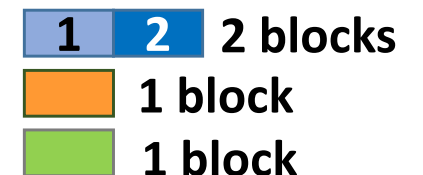


Mécanismes d'HDFS

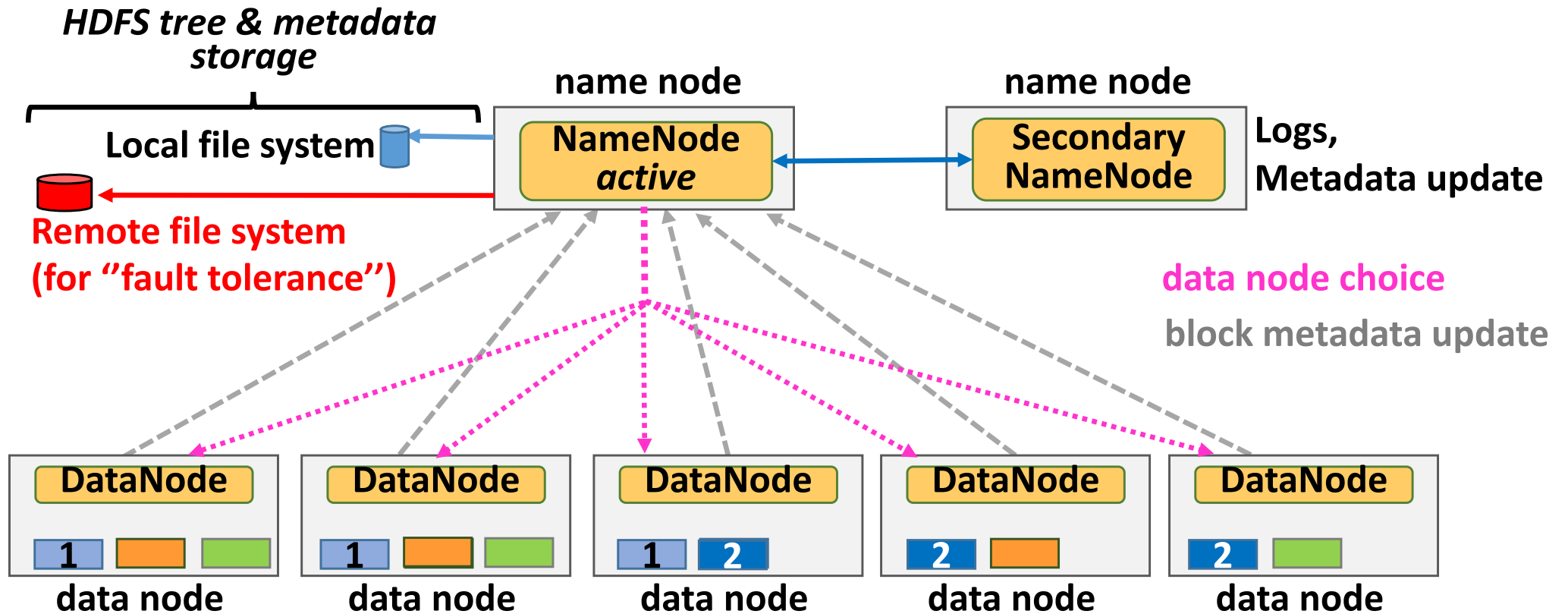


Règles de choix des nœuds :

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



Mécanismes d'HDFS



Si un nœud tombe :

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

1 2 2 blocks
 1 block
 1 block

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)

$$T_{\text{seek}} = 10\text{ms}$$

Bande passante disque std : $B_w = 100 \text{ Mo/s}$

$$T_{\text{read}} = Q / B_w$$

On veut : $T_{\text{seek}} < 1\% T_{\text{read}}$

$$\Leftrightarrow 10 \cdot 10^{-3} \text{ s} < (1/100) \cdot (Q/100) \text{ s}$$

$$\Leftrightarrow 100 \text{ (Mo)} < Q$$

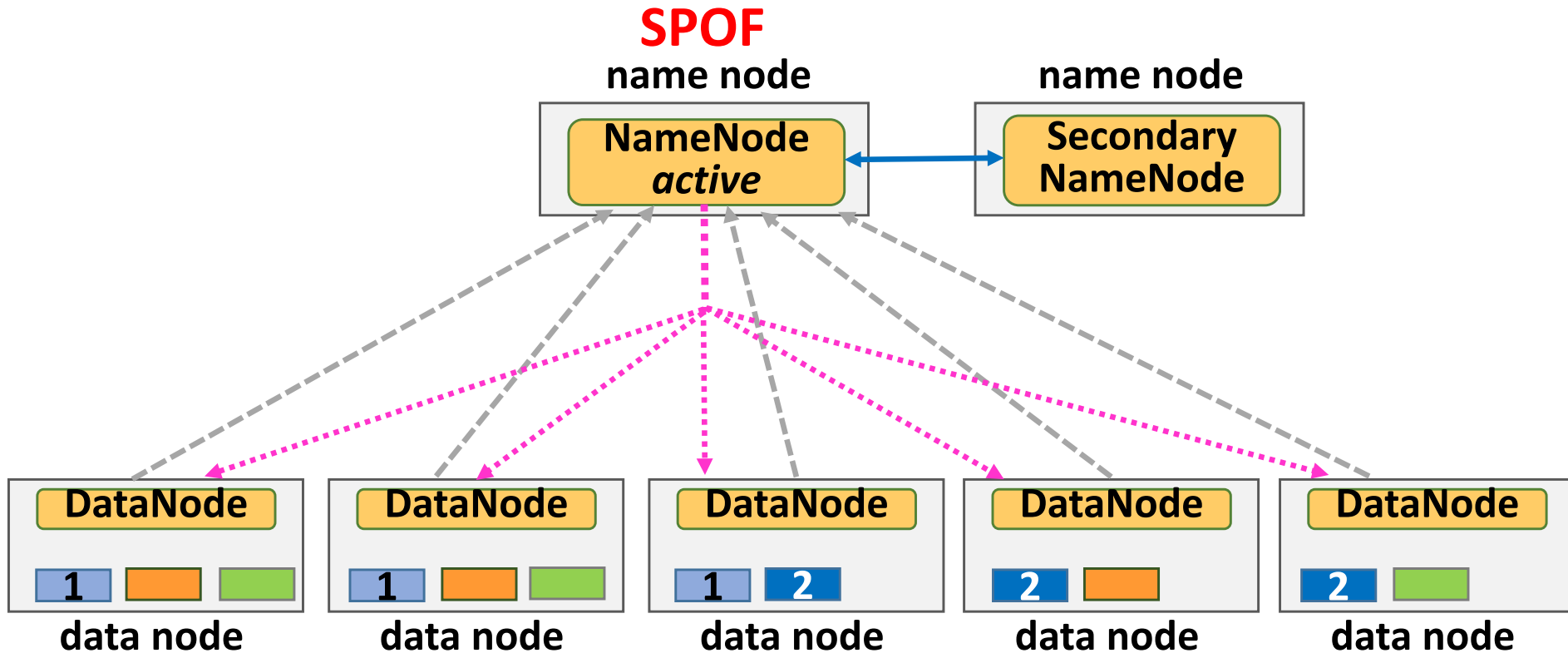
→ Des blocs de 64 Mo ou 128 Mo permettent de masquer les temps de seek

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

Principes d'HDFS

1. Framework d'Hadoop
- 2. HDFS (Hadoop Distributed File System)**
 - 2.1. Stockage des fichiers sous HDFS
 - 2.2. Mécanisme de haute disponibilité d'HDFS**
 - 2.3. Lecture de fichiers sous HDFS
 - 2.4. Ecriture de fichiers sous HDFS

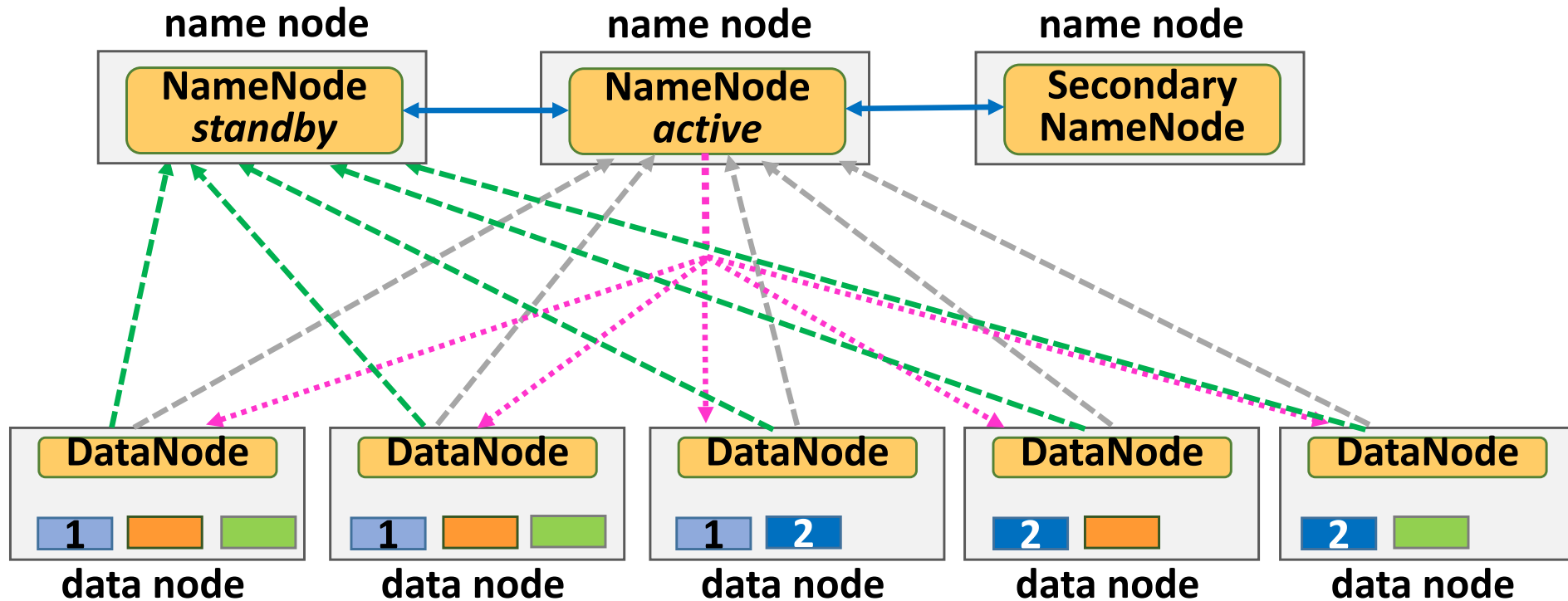
Mécanismes d'HDFS : « haute disponibilité »



Le NameNode est un SPOF : Single Point Of Failure

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

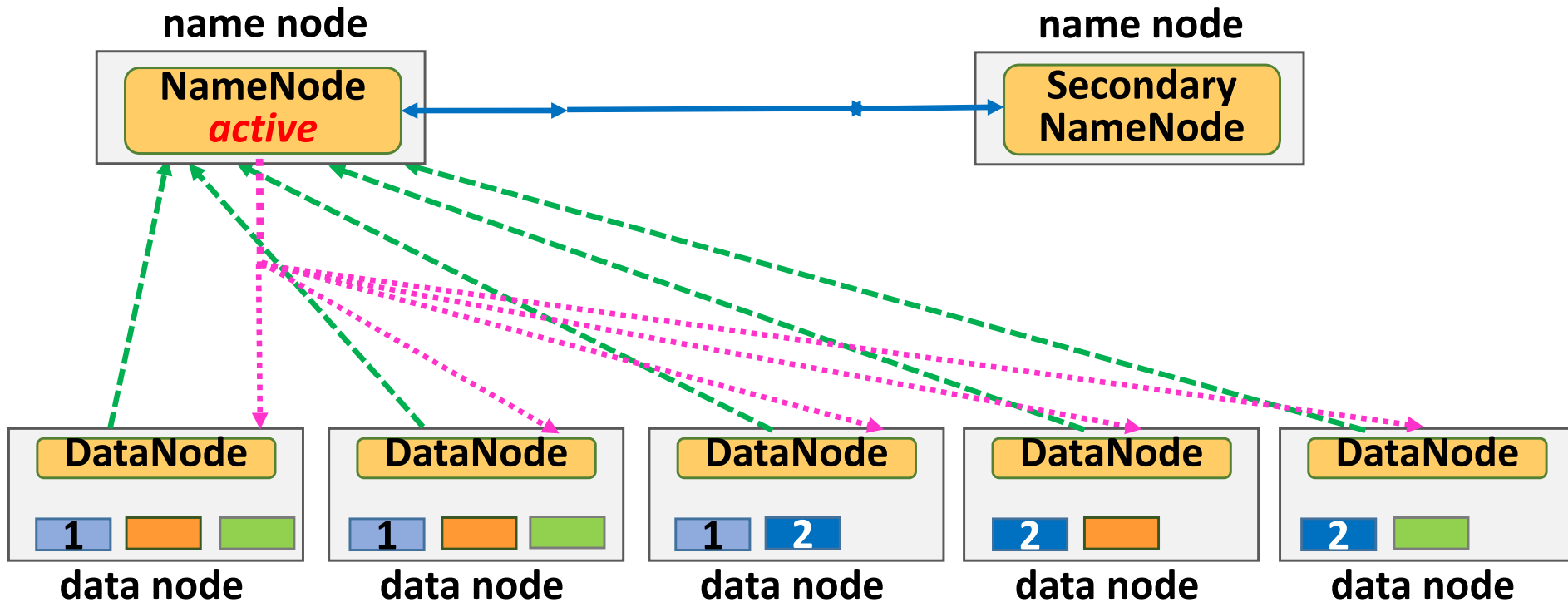
Mécanismes d'HDFS : « haute disponibilité »



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

→ On duplique le NameNode

Mécanismes d'HDFS : « haute disponibilité »

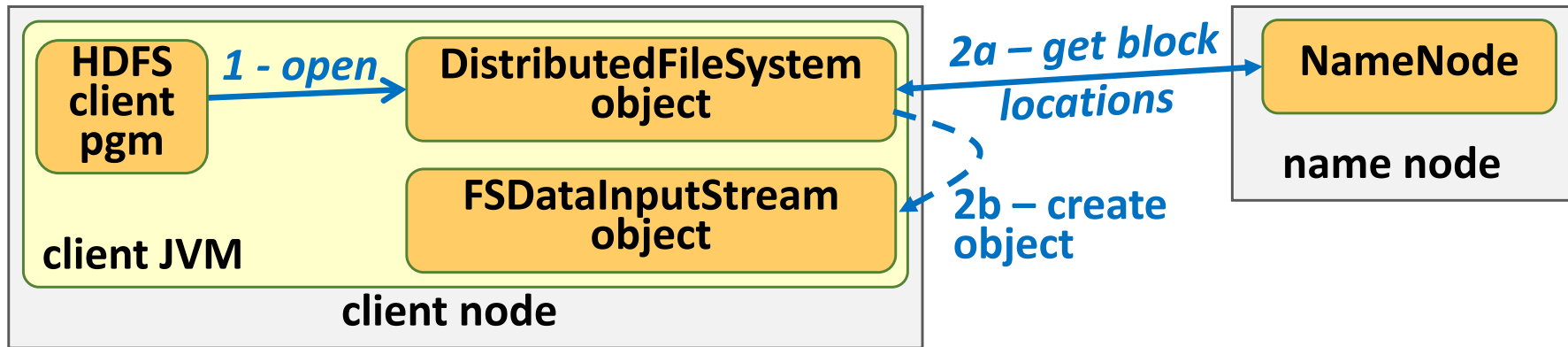


- Le NameNode en *standby* est passé *actif* « très rapidement »
 - Il n'y a plus de SPOF
- On « ne sent plus passer la panne » : Haute Disponibilité

Principes d'HDFS

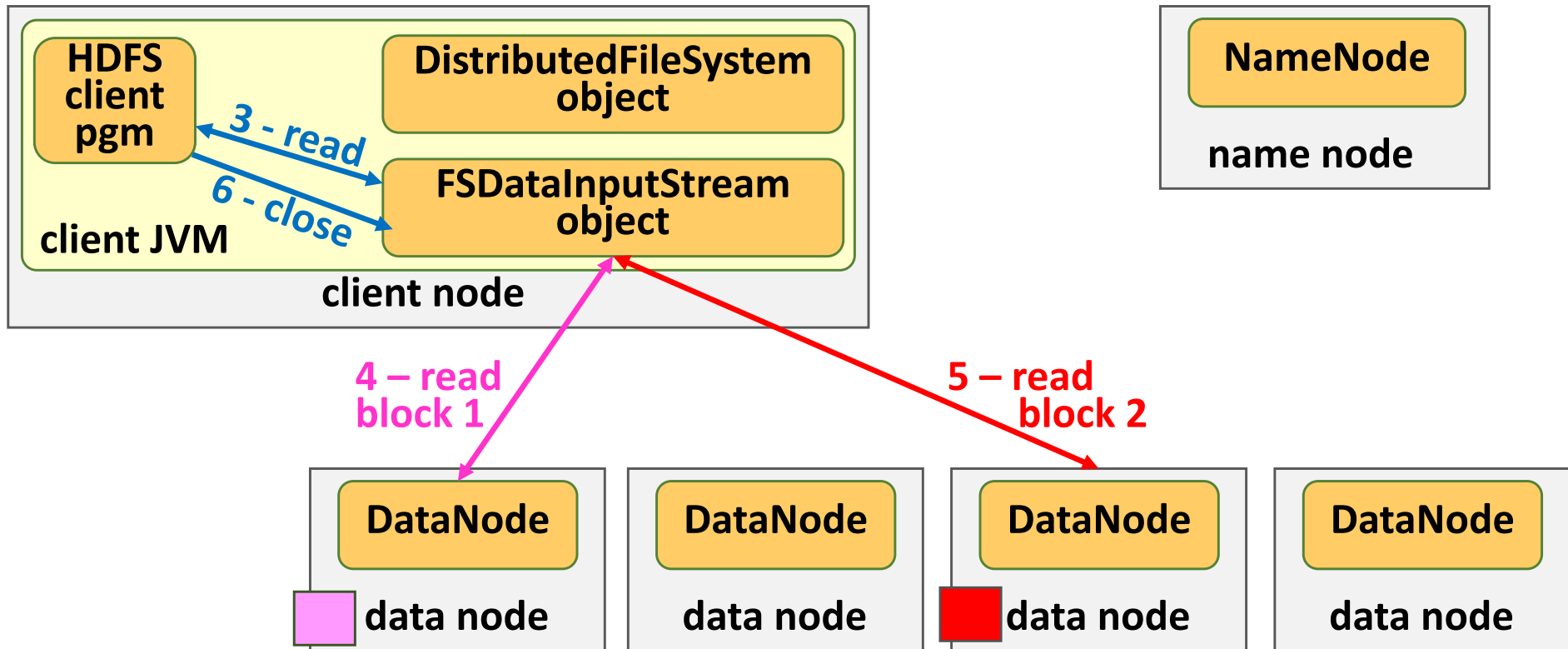
1. Framework d'Hadoop
- 2. HDFS (Hadoop Distributed File System)**
 - 2.1. Stockage des fichiers sous HDFS
 - 2.2. Mécanisme de haute disponibilité d'HDFS
 - 2.3. Lecture de fichiers sous HDFS**
 - 2.4. Ecriture de fichiers sous HDFS

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

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

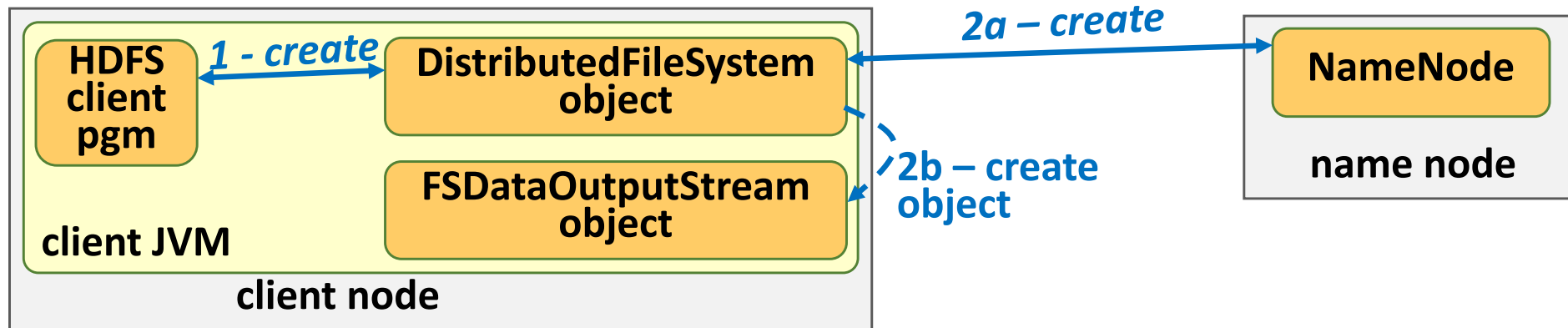
5 – Le lecteur accède au data node du second bloc et récupère ce bloc

6 – Le client referme le fichier en s'adressant au lecteur du fichier

Principes d'HDFS

1. Framework d'Hadoop
- 2. HDFS (Hadoop Distributed File System)**
 - 2.1. Stockage des fichiers sous HDFS
 - 2.2. Mécanisme de haute disponibilité d'HDFS
 - 2.3. Lecture de fichiers sous HDFS
 - 2.4. Ecriture de fichiers sous HDFS**

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 »)

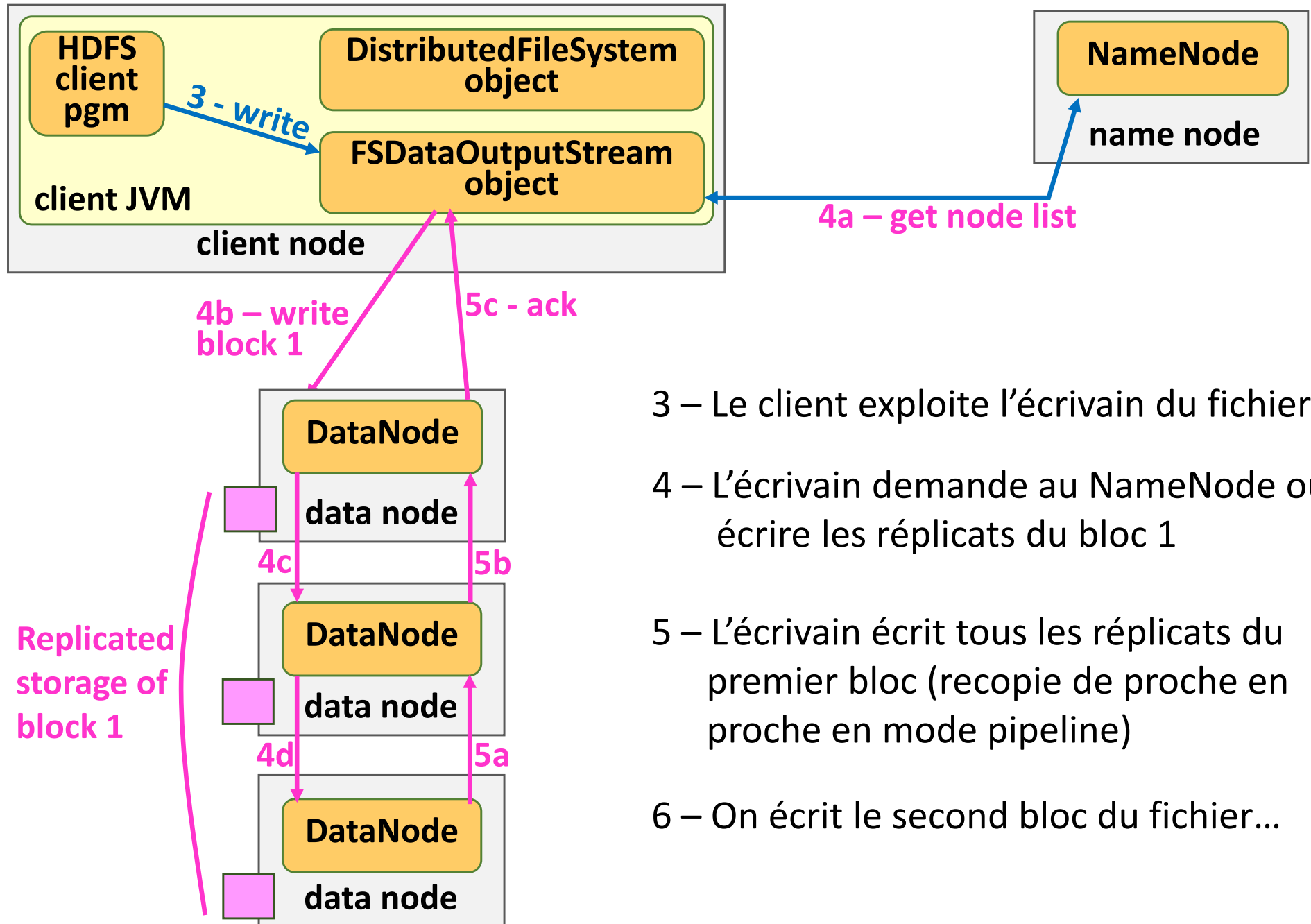
2 – Demande de localisation des nœuds d'accueil des futurs blocs du fichier

2a - Le proxy interroge le NameNode : pour savoir où écrire les blocs

Le NameNode répond au proxy

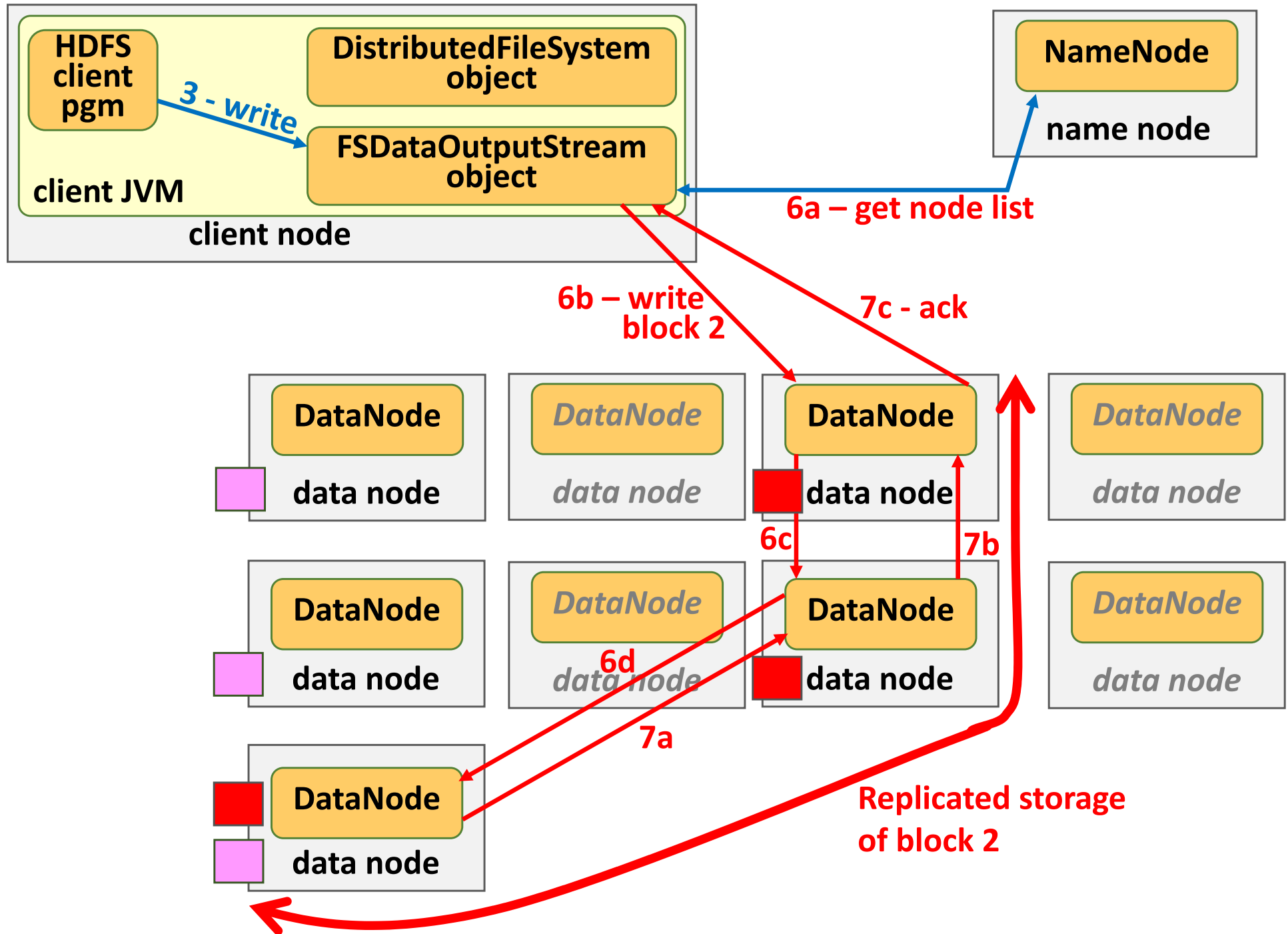
2b - Le proxy crée et retourne un objet écrivain spécialisé sur le fichier ciblé

Mécanisme d'écriture d'HDFS

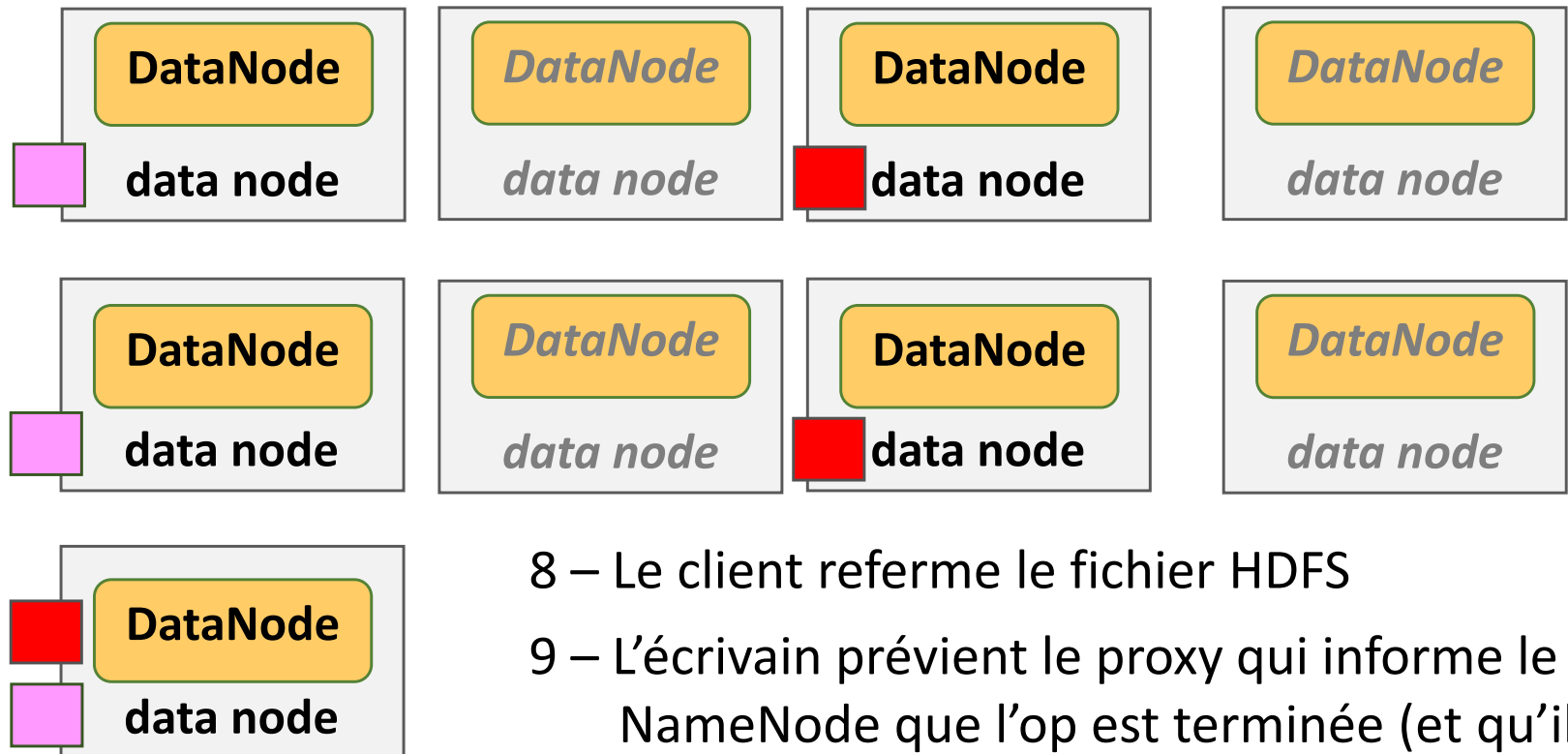
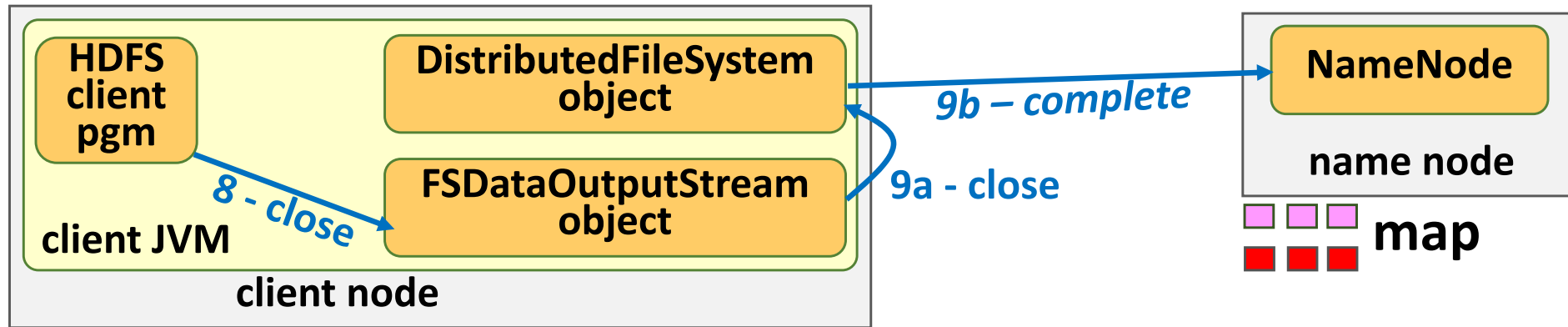


- 3 – Le client exploite l'écrivain du fichier
- 4 – L'écrivain demande au NameNode où écrire les réplicats du bloc 1
- 5 – L'écrivain écrit tous les réplicats du premier bloc (recopie de proche en proche en mode pipeline)
- 6 – On écrit le second bloc du fichier...

Mécanisme d'écriture d'HDFS



Mécanisme d'écriture d'HDFS



8 – Le client referme le fichier HDFS

9 – L'écrivain prévient le proxy qui informe le NameNode que l'op est terminée (et qu'il peut rouvrir le fichier pour un autre client)

QUIZ

Q1: 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?

QUIZ

Q2: 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?

Q3: 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?

Principes d'HDFS

