

Spark : Fonctionnement



Spark : Fonctionnement

Le principe du travail de spark est de faire du traitement de données massivement et rapidement.

Pour realiser cet objectif, spark utilise la **parrallelisation**.

Spark : Parrallelisation

La parrallelisation permet de **diviser** un problème en plusieurs sous-problèmes.

Puis de les résoudre en meme temps (parallèle) sur différents processeurs.

Le plus souvent repartie sur plusieurs machines : le cluster.

Spark : Parrallelisation

Cette **parrallelisation** est réalisée par :

- La division des données en Resilient Distributed Dataset (RDD).
 - La division du programme par le Direct Acyclic Graph.
-

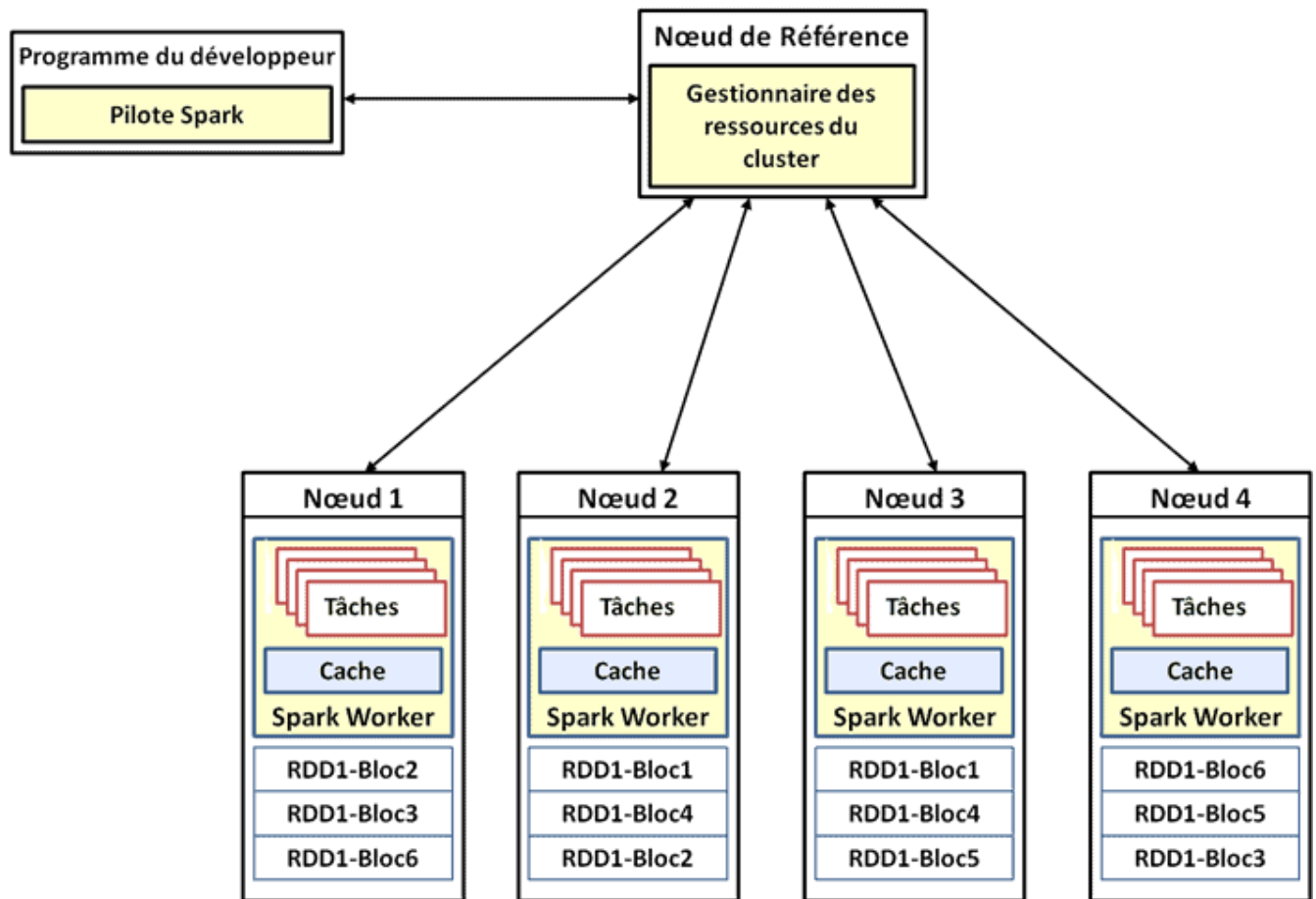
Resilient Distributed Dataset (RDD)

Pour améliorer le traitement des données, spark manipule les informations sous un format divisé en blocs d'informations réparties sur les differents noeuds.

C'est le Resilient Distributed Dataset (RDD).

Resilient Distributed Dataset (RDD)

Les blocs de données sont appelés **partitions** et sont stockés en plusieurs exemplaire sur les differents noeuds du cluster.



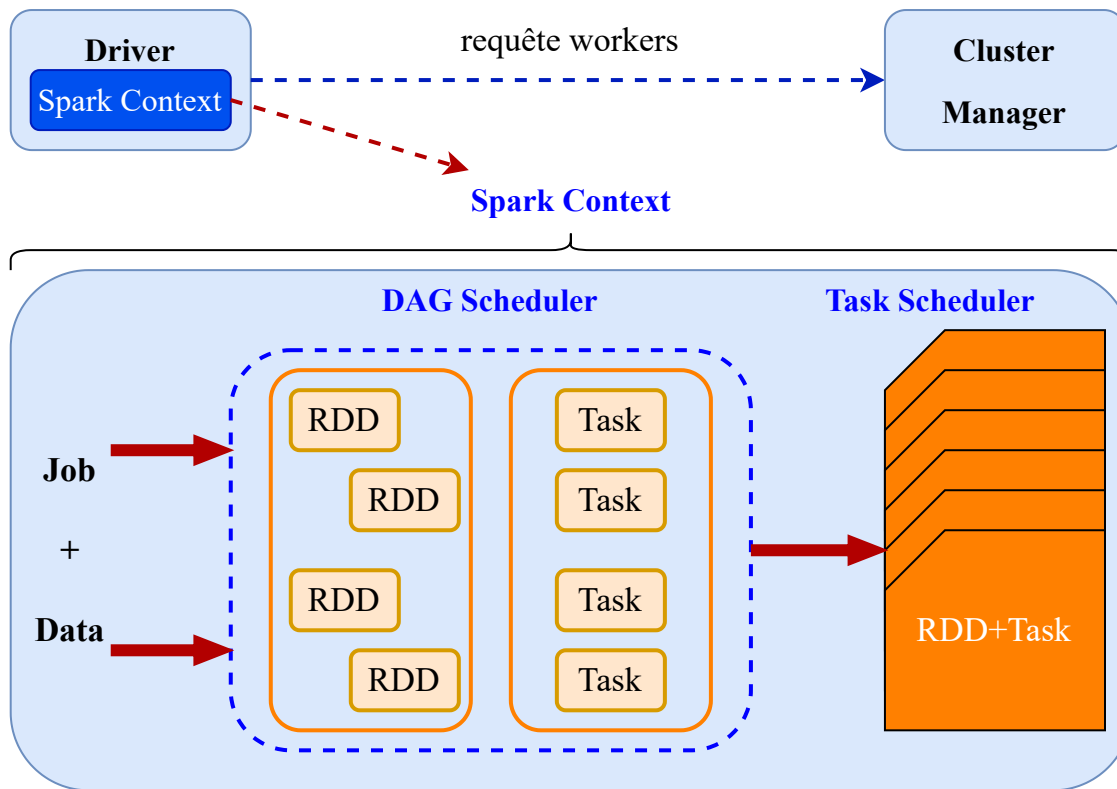
Direct Acyclic Graph (DAG)

Lorsque spark reçoit un programme, celui va être également transformé.

Le DAG va permettre de **diviser le travail** à accomplir en actions à effectuer les unes à la suite des autres.

Une tâche est ensuite planifiée pour une action sur un bloc de données.

Direct Acyclic Graph (DAG)



Resumé

Spark va diviser les données en partitions et les répartir sur les différents nœuds du cluster.

Puis va diviser le programme en tâches à effectuer sur les partitions.

Enfin, il va répartir les tâches sur les différents nœuds du cluster.

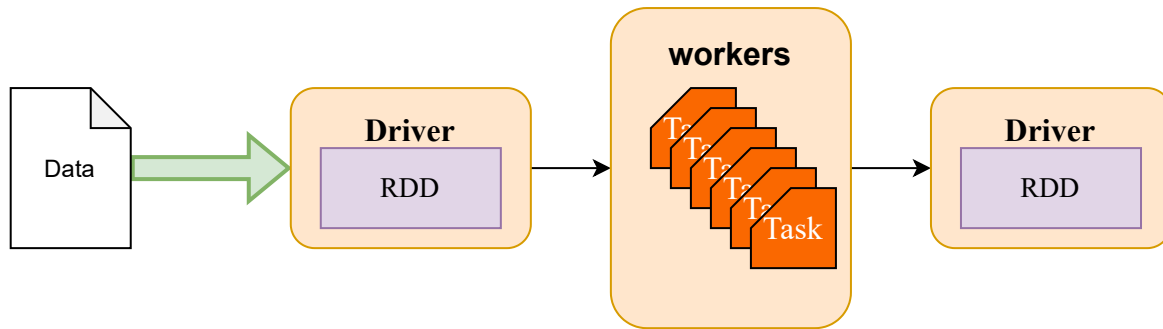
Le travail peut alors être effectué en même temps.

Spark : Flux de données

Le travail de traitement des données par Spark peut se faire selon différentes configurations:

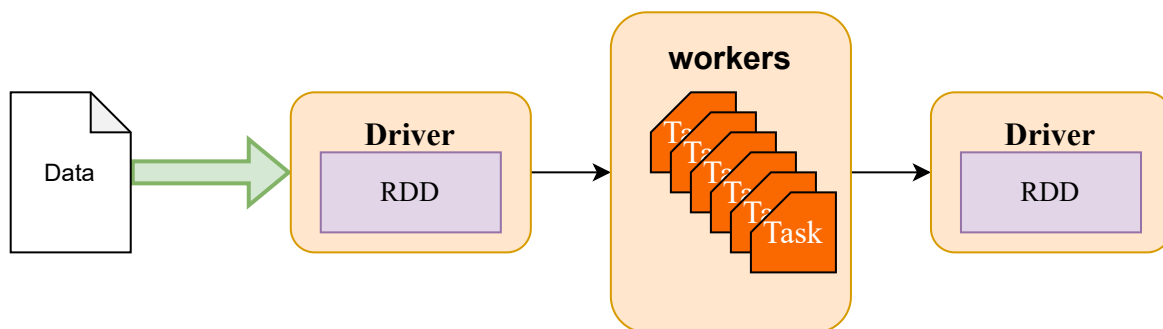
- Traitement par lots (batch)
- Traitement en micro-batch
- Traitement en temps réel (streaming)

Traitement par lots (batch)



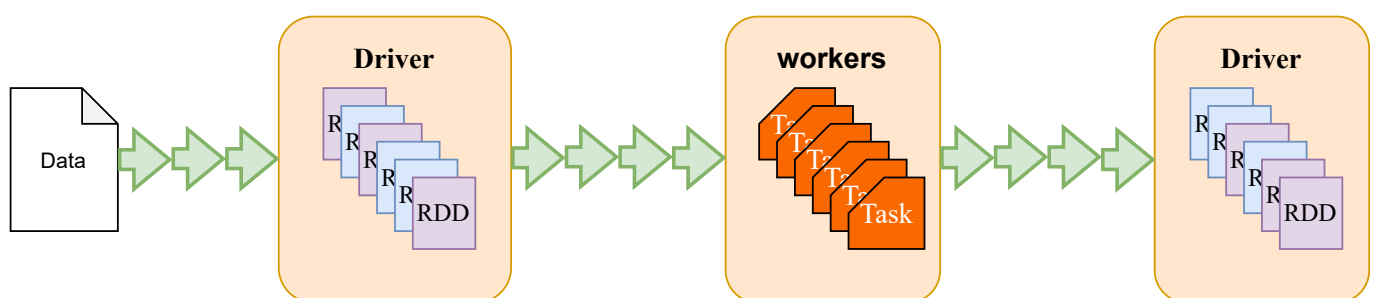
- Consiste à charger en mémoire l'ensemble des données à traiter.
- Un seule RDD sur les différents noeuds du cluster (stocké et traité de manière distribuée).

Traitement par lots (batch)



- Le traitement s'effectue une fois le chargement complet des données effectuées.
- L'ensemble est ensuite regroupé en un seul RDD.

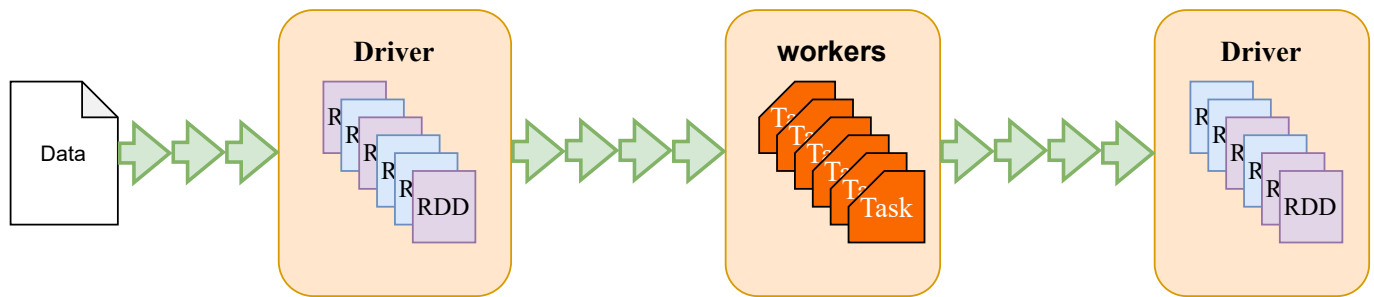
Le traitement en micro-batch



Le traitement en micro-batch consiste à charger les données par petits lots appelés **micro-batch** au lieu de tous charger d'un coup.

Interet : traitement au fur et à mesure par petit groupe sans attendre d'avoir tout récupérer.

Traitement en micro-batch



Le traitement en micro-batch est particulièrement adaptées aux données très volumineuses ou aux données incrémentales.

Traitement en temps réel (streaming)

Le traitement en temps réel consiste à traiter les données au fur et à mesure de leur arrivée.

C'est le traitement le plus rapide car les données sont traitées au fur et à mesure de leur arrivée.

Traitement en temps réel (streaming)

En pratique, le traitement en temps réel est une aggrégation des données traitées en micro-batch en continue.

[Les modules de Spark](#)