



كلية العلوم
السملاية - مراكش
FACULTÉ DES SCIENCES
SEMLALIA - MARRAKECH

Université Cadi Ayyad

Faculté des Sciences Semlalia
Département d'Informatique

Projet de Big Data Labs

Pour la validation de module en Master Science de données.

L'analyse des films regardés sur PoliTV.

Réalisé par :

Omar BOUDELLAH && Mustapha AHRICHA

Encadré par :

Pr. EL ALAOUI Hasna

Professeur, Département informatique, FSSM Marrakech

Année Universitaire : 2023 / 2024

Table de matière

Chapitre 1: Introduction	3
1- Objectifs du projet :	3
2- Importance du projet dans le contexte du Big Data	4
Chapitre 2: Cadre théorique et technologique	4
1- Big data	4
1.0 Revue de la littérature sur les technologies Big Data	4
1.1 Introduction au Big Data	4
1.2 Caractéristiques du Big Data	5
1.3 Importance du Big Data	6
1.4 Exemple use case d'un projet big data	6
1.5 Sources de Big Data	6
1.6 Technologies de Big Data	7
2- Hadoop	7
2.1 Introduction	7
2.2 Hadoop 1	9
2.3 Hadoop 2	12
3- Présentation de Hadoop MapReduce dans le cas de notre projet	14
3.1 Introduction à Hadoop MapReduce	14
3.2 Fonctionnement de MapReduce	14
3.3 Écosystème Hadoop	15
3. 4. Avantages de Hadoop MapReduce	15
Conclusion	16
4- Présentation de Spark et RDD	16
4.1 Introduction à Apache Spark et RDD	16
4.2 Concept de RDD (Resilient Distributed Datasets)	17
4.3 Avantages de Spark et RDD	18
Conclusion	19
Chapitre 3: Cas pratique de Projet	19
1- Les données utilisées	19
2- Code source Hadoop MapReduce	20
=> Mapper class	21
=> Reduce class	21
Résultat:	22
3- Code source Spark et RDD	22
Question 1:	22
Question2:	23

Table de figures

Figure 1	Big data	5
Figure 2	5 V de big data	5
Figure 3	Big data au sein des unités de production	6
Figure 4	Logo Hadoop	7
Figure 5	Logo Spark	7
Figure 6	Composants fondamentaux	9
Figure 7	Principe de base	10
Figure 8	HDFS	10
Figure 9	Exemple d'illustration	11
Figure 10	schema MapReduce1	11
Figure 11	Architecture hadoop 2	12
Figure 12	MapReduce1 Vs MapReduce2	13
Figure 13	Hadoop1 & Hadoop2	13
Figure 14	Use case MapReduce	14
Figure 15	Opérations sur RDD	18
Figure 16	Bibliothèque Etendues	18
Figure 17	Fichiers de travail	20
Figure 18	Code source Hadoop MapReduce	21
Figure 19	Mapper Class	21
Figure 20	Reduce class	22
Figure 21	Resultat partie 1	22
Figure 22	Code question1 spark	23
Figure 23	Output question1 spark	23
Figure 24	Code question2 spark	24
Figure 25	Output question2 spark	24

Chapitre 1: Introduction

1-Objectifs du projet :

Le projet vise à répondre aux besoins d'analyse spécifiques des responsables de PoliTV sur les films regardés. Les principales tâches à accomplir dans le cadre de ce projet consistent à concevoir une application unique, en utilisant à la fois Hadoop MapReduce (pour le premier point) et Spark avec RDD (pour le deuxième point), et à écrire le code Java et Python correspondant. Les objectifs spécifiques sont les suivants :

1. Films Regardés par un Seul Utilisateur au Cours de l'Année 2019 (Hadoop MapReduce).
2. Films Regardés Fréquemment mais Seulement Pendant un An au Cours des Cinq Dernières Années (Spark et RDD) .
- 3.Extraction des Films les Plus Populaires Depuis au Moins Deux Ans (Spark et RDD) .

2-Importance du projet dans le contexte du Big Data

les analyses que nous faisons permettront aux responsables de PoliTV de mieux comprendre les habitudes de visionnage de leurs utilisateurs, d'identifier les films populaires sur une période prolongée, et d'optimiser leurs recommandations de contenu. En répondant à ces exigences, notre application contribuera à l'amélioration de la plateforme PoliTV en fournissant des informations précieuses sur les préférences des utilisateurs et la popularité des films au fil du temps.

Chapitre 2: Cadre théorique et technologique

1- Big data

1.0 Revue de la littérature sur les technologies Big Data

La revue de la littérature sur les technologies Big Data est une section cruciale qui permet d'explorer les concepts, les défis, et les évolutions significatives dans le domaine du Big Data. Voici un exemple de contenu que vous pourriez inclure dans cette section :

1.1 Introduction au Big Data

Le Big Data, terme émergé au cours des dernières années, représente un changement fondamental dans la manière dont nous collectons, stockons, traitons et analysons les données. Il est caractérisé par le volume massif, la variété de formats, la vitesse de génération et la nécessité d'une vélocité de traitement rapide.

➡ **Définition Big Data** : Traitement rapide de grands volume de données difficilement gérables avec les outils traditionnels (une base de données classique ou un seul ordinateur).

Le terme Big data focalise sur des données qui contient une plus grand variété, arrivant dans des volumes croissants et une plus de vitesse.

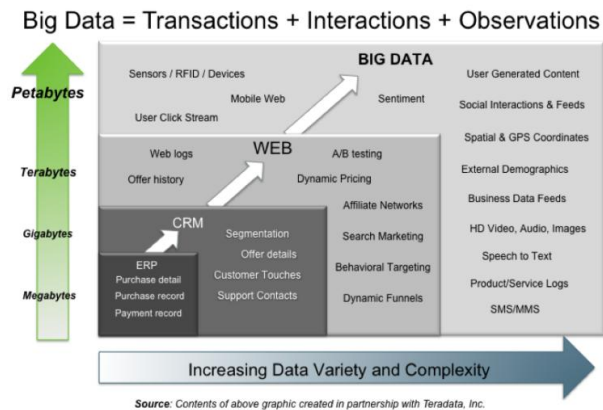


Figure 1 Big data

1.2 Caractéristiques du Big Data

1. **Volume** : Les données sont générées à une échelle exponentielle, dépassant les capacités des systèmes traditionnels de gestion de bases de données.
2. **Variété** : Les données proviennent de diverses sources telles que les réseaux sociaux, les capteurs, les transactions en ligne, etc., et peuvent être structurées, semi-structurées ou non structurées.
3. **Vélocité** : Les données doivent être traitées en temps réel ou quasi-réel, car elles sont générées rapidement.
4. **Véracité** : En raison de la diversité des sources, la qualité des données peut varier, nécessitant des techniques avancées de vérification.
5. **Valeur** : Au-delà de simplement collecter et stocker les données, il est crucial d'extraire des informations significatives et exploitables. La capacité à analyser ces données pour en extraire de la valeur, que ce soit pour prendre des décisions commerciales, améliorer les processus, ou comprendre les modèles, est essentielle. Cette valeur peut être obtenue par le biais d'analyses avancées, de l'apprentissage automatique ou d'autres techniques d'analyse de données.

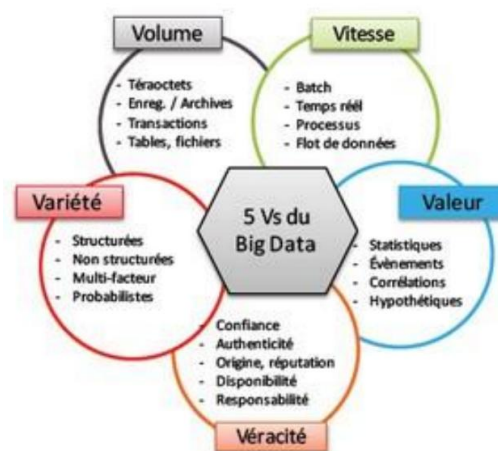


Figure 2 5 V de big data

1.3 Importance du Big Data

L'émergence du Big Data a un impact significatif dans divers domaines :

1. **Recherche Scientifique** : Facilite la gestion et l'analyse de grandes quantités de données expérimentales.
2. **Entreprises** : Permet une prise de décision basée sur les données, l'optimisation des processus et la personnalisation des services.
3. **Santé** : Facilite la recherche médicale, la gestion des dossiers médicaux et la surveillance des épidémies.
4. **Gouvernance** : Aide à l'analyse des tendances socio-économiques, à la prévention de la criminalité, etc.

1.4 Exemple use case d'un projet big data

Exemple - Usine 4.0

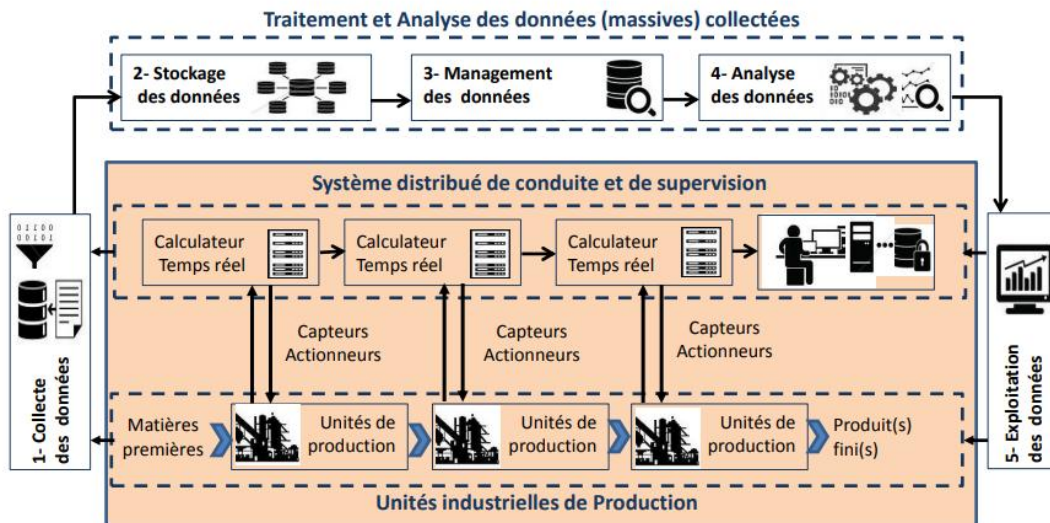


Figure 3 Big data au sein des unités de production

1.5 Sources de Big Data

Les principales sources de données comprennent :

1. Données des Médias Sociaux : Tweets, publications sur Facebook, vidéos sur YouTube, etc.
2. Capteurs et IoT (Internet des Objets) : Données générées par des appareils connectés.
3. Transactions en Ligne : Historique d'achats, transactions financières, etc.
4. Données Générées par l'Utilisateur : Recherches en ligne, commentaires, avis, etc.

1.6 Technologies de Big Data

Les technologies Big Data englobent un éventail d'outils et de techniques utilisés pour collecter, stocker, traiter et analyser ces énormes quantités de données.

Les principaux technologies clés associées au Big Data :

- Hadoop: Un framework open source qui permet de stocker et de traiter de gros volumes de données sur des clusters de serveurs. Il utilise le modèle de programmation MapReduce pour le traitement distribué.



Figure 4 Logo Hadoop

- Spark: Une plateforme de traitement de données rapide et polyvalente qui fonctionne en mémoire et peut traiter des charges de travail plus diverses que Hadoop. Il est souvent utilisé pour le traitement en temps réel et les analyses avancées.



Figure 5 Logo Spark

2- Hadoop

2.1 Introduction

Hadoop a évolué pour devenir l'un des piliers du traitement et du stockage de données massives à l'ère du Big Data. À ses débuts, Hadoop était principalement axé sur deux composants clés : le stockage et le traitement. Alors c'est quoi hadoop? Et quel est son role?

2.1.1 Définition

Hadoop est un framework open source conçu pour le stockage et le traitement distribué de grands ensembles de données sur des clusters d'ordinateurs. Il est principalement utilisé pour gérer et analyser de grandes quantités de données structurées ou non structurées, offrant ainsi une solution évolutive pour le traitement de données massives.

Il est conçu par Doug Cutting en 2004, il est introduit par Google en s'inspirant des solutions Google comme MapReduce, GoogleFS et BigTable.

Déployable et configurable très facilement.

2.1.2 Utilisabilité

L'utilisabilité de Hadoop dépend des besoins spécifiques de chaque projet :

1. **Traitement de gros volumes de données** : Hadoop est particulièrement efficace pour le traitement de grandes quantités de données réparties sur plusieurs nœuds dans un cluster. Si votre cas d'utilisation implique des volumes massifs de données, Hadoop peut être utile.
2. **Évolutivité** : Il offre une grande évolutivité horizontale, ce qui signifie que vous pouvez ajouter de nouveaux nœuds au cluster pour augmenter la capacité de stockage et de traitement.
3. **Tolérance aux pannes** : Hadoop est conçu pour gérer les pannes matérielles. Il réplique automatiquement les données sur plusieurs nœuds, assurant ainsi une certaine tolérance aux pannes.
4. **Adaptabilité** : Il est assez souple pour s'intégrer à d'autres outils et technologies. Cependant, sa configuration et sa gestion peuvent être complexes.
5. **Complexité** : La mise en place et la configuration initiales peuvent être complexes, nécessitant des compétences spécifiques en administration de systèmes distribués.
6. **Performance grâce au parallélisme des données** : Hadoop divise les tâches en petits morceaux et les exécute simultanément sur différents nœuds, exploitant ainsi le parallélisme pour accélérer le traitement des données massives.

2.1.3 Composants fondamentaux

Hadoop est constitué de 2 grandes parties :

HDFS (Hadoop Distributed File Système) : destiné pour le stockage distribué des données.

MapReduce (Distributed Programming Framework) : destine pour le traitement distribué des données.

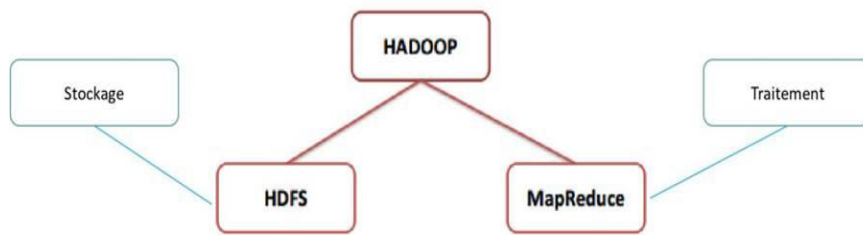


Figure 6 Composants fondamentaux

➡ Version Hadoop

Il existe deux principales versions de Hadoop : Hadoop 1 et Hadoop 2.

2.2 Hadoop 1

2.2.1 Introduction

Hadoop 1 était la première version majeure du framework. Il comprenait deux composants principaux :

- le système de stockage Hadoop Distributed File System (HDFS) et le framework de traitement MapReduce.
- MapReduce était utilisé pour traiter et analyser les données stockées dans HDFS. Cependant, cette version avait quelques limitations, notamment en termes d'extensibilité et de types de traitement supportés.

2.2.2 Principe de base

Le principe de base de Hadoop 1 repose sur deux composants principaux : Hadoop Distributed File System (HDFS) et le framework de traitement MapReduce.

- ❖ Copier les données vers HDFS: les données sont divisées et stockées sur un ensemble de nœuds.
- ❖ Traiter les données là où elles sont stockées (MAP) et recueillir les résultats (REDUCE).
- ❖ Copier les résultats vers HDFS.

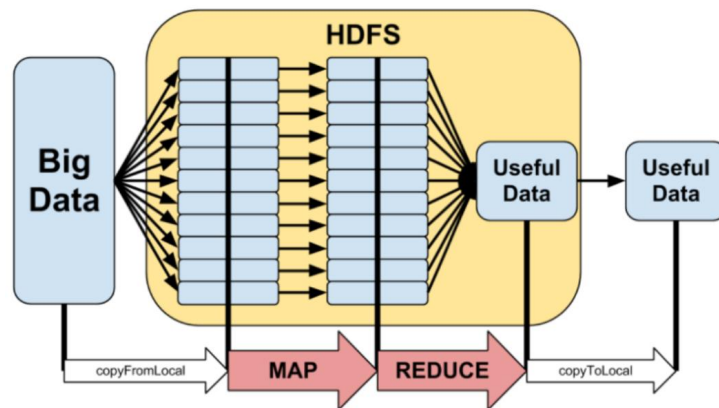


Figure 7 Principe de base

2.2.3 Noyaux de Hadoop

Hadoop Distributed File System (HDFS) :

HDFS est conçu pour stocker de très gros fichiers de données en les répartissant sur plusieurs nœuds dans un cluster. Il divise les fichiers en blocs et les réplique sur différents nœuds pour assurer la tolérance aux pannes et la disponibilité des données.

Il fonctionne selon un modèle maître-esclave, où un nœud maître, le NameNode, gère le système de fichiers en conservant les métadonnées, tandis que les nœuds esclaves, les DataNodes, stockent les données réelles.

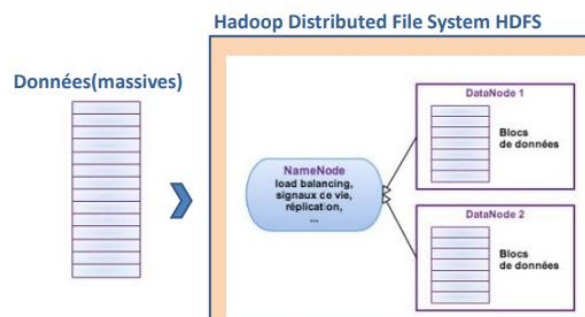


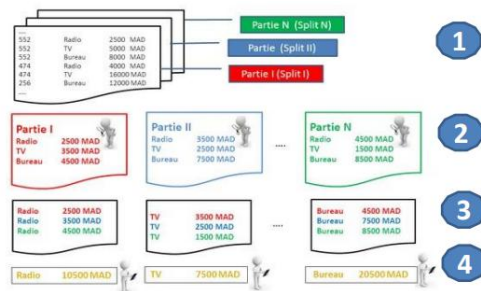
Figure 8 HDFS

MapReduce :

MapReduce est un modèle de programmation pour le traitement parallèle et distribué de grands ensembles de données. Il divise les tâches en deux phases principales : la phase de Map où les données sont traitées et filtrées, puis la phase de Reduce où les résultats intermédiaires sont combinés pour produire le résultat final.

MapReduce fonctionne en distribuant ces tâches sur les nœuds du cluster, ce qui permet un traitement parallèle et efficace des données stockées dans HDFS.

➡ MapReduce: Exemple d'illustration



1. Les données initiales sont scindées en lots (**split**)
2. Chaque lot est confié à un *Mapper* qui évalue la **fonction Map** sur chaque enregistrement
3. Les listes obtenues en sortie des *Mappers* sont concaténées, classées par valeur de la clé et scindées en lots par le Framework. C'est la phase dite du **Shuffle**.
4. La **fonction Reduce** exécutée par chaque *Reducer* agrège toutes les valeurs intermédiaires associées à une même clé.



Figure 9 Exemple d'illustration

2.2.4 Infrastructure(s) logicielle(s) d'implémentation du pattern MapReduce

Job MapReduce: Application client qui s'exécute dans Hadoop

- Chaque Job est découpé en tâches **Map et/ou Reduce**.
- **JobTracker** est un processus maître dont le rôle est de planifier l'exécution des tâches (Map ou Reduce) qu'il attribue aux processus esclaves appelés TaskTracker.
- Chaque **TaskTracker** doit envoyer des notifications régulières, appelées Heartbeat Call, pour informer le JobTracker de l'état d'avancement des tâches qui lui ont été confiées.
- Le **JobTracker** se charge à la fois: allouer des ressources (CPU, mémoire, ... etc.) aux différentes tâches, coordonner les Jobs MapReduce, gérer les fautes en réallouant les slots éventuellement, ... etc.

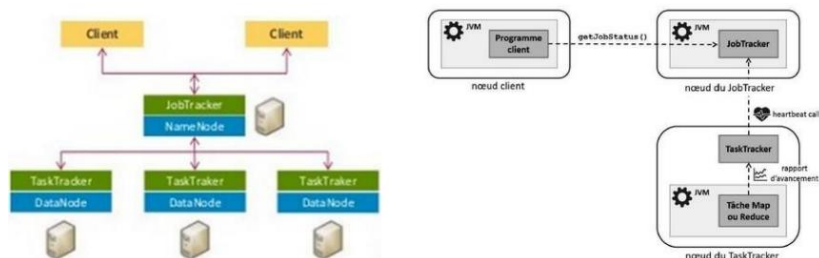


Figure 10 schema MapReduce1

2.3 Hadoop 2

2.3.1 Introduction

Hadoop 2.x a introduit plusieurs améliorations majeures par rapport à la première version. La modification la plus significative était l'ajout du framework **YARN** (Yet Another Resource Negotiator).

- **YARN** (qui signifie "Yet Another Resource Negotiator") a introduit une couche de gestion des ressources plus flexible dans l'écosystème Hadoop, permettant la prise en charge de différents types de frameworks de traitement, au-delà de MapReduce.

Cela a ouvert la porte à d'autres modèles de calcul, comme Apache Spark, Apache Tez, et d'autres, offrant ainsi une plus grande flexibilité dans le traitement des données et une meilleure utilisation des ressources.

2.3.2 Architecture de Hadoop 2

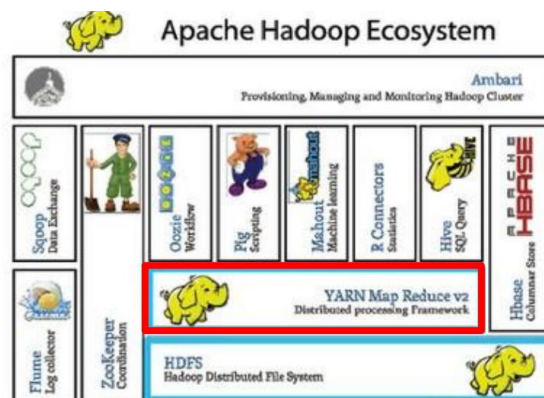


Figure 11 Architecture hadoop 2

YARN signifie Yet Another Resource Negotiator. C'est une plateforme générique de gestion des ressources dans un cluster. YARN a été introduit avec Hadoop 2.0, un framework de traitement distribué open source d'Apache.

YARN sépare le **JobTracker** (Hadoop 1.x) en un **ResourceManager** (RM) au niveau du cluster et un **ApplicationMaster** (AM) spécifique à l'application.

L'architecture de YARN suit un modèle architectural maître-esclave : **ResourceManager** en tant que maître et **NodeManager** spécifique à chaque nœud en tant qu'esclave. NodeManager remplace le **TaskTracker** dans Hadoop 1.x.

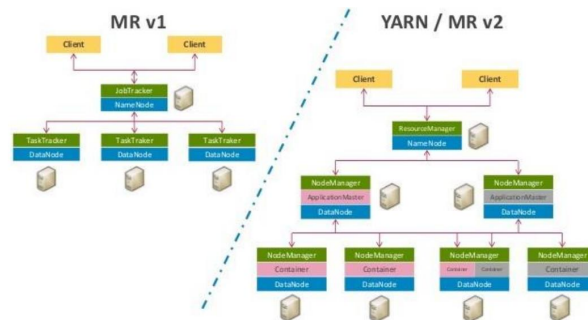


Figure 12 MapReduce1 Vs MapReduce2

2.3.3 Différence Hadoop 1 VS Hadoop 2

En tant que système d'exploitation de données pour Hadoop 2.0, YARN permet à plusieurs applications de coexister dans le même cluster partagé. YARN prend également en compte la scalabilité, une utilisation élevée, prend en charge plusieurs modèles de programmation, propose un modèle de ressources flexible, la sécurité, la fiabilité et la compatibilité.

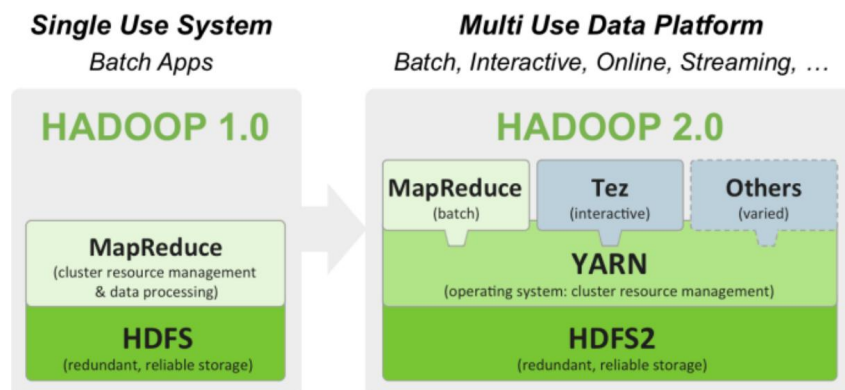
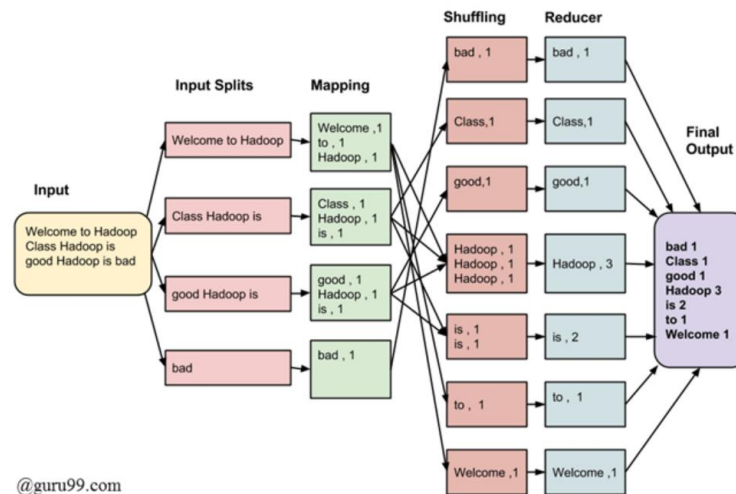


Figure 13 Hadoop1 & Hadoop2

3- Présentation de Hadoop MapReduce dans le cas de notre projet



@guru99.com

Architecture MapReduce

Figure 14 Use case MapReduce

3.1 Introduction à Hadoop MapReduce

3.1.1 Contexte et Origines

Hadoop MapReduce est un modèle de programmation et un framework open-source développé par Apache pour le traitement distribué de données volumineuses. Il trouve ses racines dans les travaux de Google sur le traitement distribué de données, en particulier le modèle MapReduce.

3.1.2 Objectif de MapReduce

L'objectif principal de MapReduce est de fournir une approche scalable et parallèle pour traiter de grandes quantités de données sur des clusters de machines.

3.2 Fonctionnement de MapReduce

3.2.1 Phases de MapReduce

MapReduce fonctionne en deux phases principales : la phase de Map et la phase de Reduce.

- Phase de Map : Dans cette phase, les données sont divisées en blocs, puis traitées en parallèle par plusieurs mappers. Chaque mapper produit une sortie intermédiaire constituée de paires clé-valeur.

- Phase de Reduce : Les résultats intermédiaires sont triés et regroupés par clé, puis traités par plusieurs reducers qui effectuent une opération de réduction pour générer le résultat final.

3.2.2 Étapes du Processus

- Input Splitting : Les données sont divisées en blocs gérables appelés "splits".
- Mapping : Application de la fonction map à chaque split, générant des paires clé-valeur intermédiaires.
- Shuffling : Tri et regroupement des résultats intermédiaires par clé.
- Reducing : Application de la fonction reduce à chaque groupe de données associées à une clé, produisant le résultat final.

3.3 Écosystème Hadoop

3.3.1 HDFS (Hadoop Distributed File System)

HDFS est le système de fichiers distribué utilisé par Hadoop pour le stockage des données. Il divise les données en blocs et les distribue sur plusieurs nœuds du cluster.

3.3.2 YARN (Yet Another Resource Negotiator)

YARN est le gestionnaire de ressources qui gère et planifie les ressources du cluster, permettant l'exécution de diverses applications, y compris MapReduce, de manière concurrente.

3. 4. Avantages de Hadoop MapReduce

3.4.1 Scalabilité

MapReduce offre une scalabilité horizontale, permettant d'ajouter simplement des nœuds au cluster pour traiter davantage de données.

3.4.2 Tolérance aux Pannes

Le modèle MapReduce est conçu pour être tolérant aux pannes, avec la capacité de récupérer automatiquement à partir d'échecs.

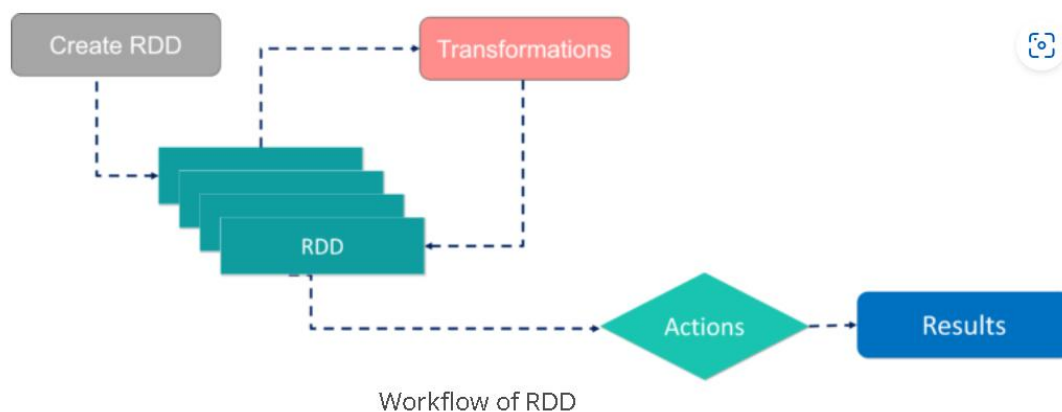
3.4.3 Facilité d'Utilisation

Hadoop MapReduce simplifie le développement en fournissant un modèle de programmation abstrait pour les développeurs.

Conclusion

En conclusion, Hadoop MapReduce est un framework puissant pour le traitement distribué de données massives, offrant une approche scalable et tolérante aux pannes. Bien que son modèle ait été révolutionnaire, l'évolution continue du domaine du Big Data a conduit au développement de frameworks plus modernes, tels que Apache Spark, pour relever de nouveaux défis.

4- Présentation de Spark et RDD



4.1 Introduction à Apache Spark et RDD

4.1.1 Contexte et Origines

Apache Spark est un framework open-source de traitement de données en temps réel et par lots,

émergeant comme une évolution significative dans le domaine du Big Data. Conçu pour surmonter les limitations de MapReduce, Spark a rapidement gagné en popularité grâce à ses performances améliorées et à sa polyvalence.

4.1.2 Objectif de Spark

Spark vise à offrir une solution plus rapide et flexible pour le traitement distribué des données, en introduisant des améliorations majeures par rapport au modèle MapReduce. Il a été développé pour répondre aux besoins croissants de traitement en temps réel et de diverses applications d'analyse de données.

4.2 Concept de RDD (Resilient Distributed Datasets)

4.2.1 Définition de RDD

Les RDD, ou Resilient Distributed Datasets, constituent le cœur de Spark. Un RDD représente une collection immuable, partitionnée et distribuée d'objets sur un cluster. L'immuabilité garantit la résilience, et la distribution favorise le traitement parallèle sur plusieurs nœuds.

4.2.2 Caractéristiques de RDD

- Résilience : Les RDD sont résilients, ce qui signifie qu'ils peuvent être reconstitués en cas de défaillance d'un nœud du cluster.
- Distribution : Les données sont distribuées de manière transparente, permettant un traitement parallèle efficace.

4.2.3 Opérations sur RDD

- Transformation : Les transformations créent de nouveaux RDD à partir d'un existant. Par exemple, map, filter, reduceByKey.
- Action : Les actions déclenchent le calcul sur le RDD, renvoyant les résultats au programme. Exemples : count, collect, saveAsTextFile.

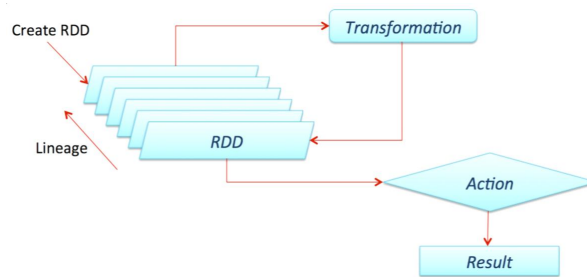


Figure 15 Opérations sur RDD

4.3 Avantages de Spark et RDD

4.3.1 Performances Améliorées

Spark offre des performances supérieures à MapReduce en minimisant les E/S disque et en favorisant le traitement en mémoire, ce qui conduit à des temps d'exécution plus courts.

4.3.2 Prise en Charge du Traitement en Temps Réel

Avec Spark Streaming, Spark excelle dans le traitement en temps réel, permettant des analyses continues sur les flux de données.

4.3.3 Bibliothèques Étendues

Spark propose un écosystème riche de bibliothèques, telles que Spark SQL pour le traitement SQL, MLlib pour l'apprentissage automatique, et GraphX pour le traitement des graphes.

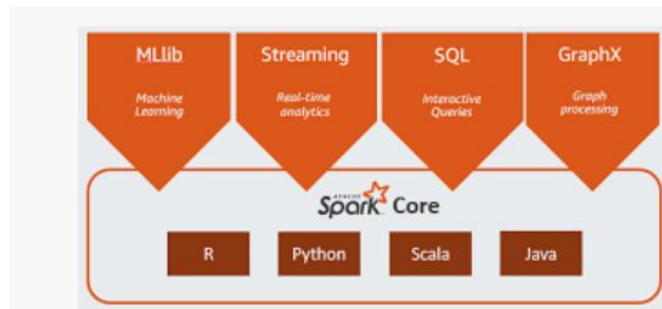


Figure 16 Bibliothèque Étendue

Conclusion

En conclusion, Apache Spark, avec son modèle RDD, représente une avancée significative dans le traitement distribué des données. Ses caractéristiques de résilience, de distribution et de performances en font un choix privilégié dans le domaine du Big Data, y compris notre projet axé sur l'analyse des films regardés.

Chapitre 3: Cas pratique de Projet

1- Les données utilisées

Pour les données, nous utilisons les données fournies dans la Classroom. Cependant, nous considérons que le fichier qui est fréquemment regardé par des utilisateurs **3 fois au lieu de 1000**, en réponse à la question 1 de Spark. De plus, nous supprimons une ligne dans le fichier 'watchedMovies' pour répondre à la question 2 de la partie 2. Et voilà le contenu de tous les fichiers de projet:

Fichier User.txt

```
PaoloG76, Male, 1976, Italy  
User1, Male, 1976, Italy  
User2, Female, 1986, France
```

Fichier Movies.txt

```
MID124, Ghostbusters, Ivan Reitman, 1984/05/01  
MID12, Title12, Ivan Reitman, 1984/05/01  
MID10, Title10, Ivan Reitman, 1984/05/01
```

Fichier WatcheMovies.txt

```

PaoloG76,MID124,2018/06/01_14:18,2018/06/01_16:10
PaoloG76,MID124,2018/07/01_14:18,2018/07/01_16:10
User1,MID124,2018/07/01_14:18,2018/07/01_16:10
User2,MID124,2018/08/01_14:18,2018/08/01_16:10
PaoloG76,MID12,2018/06/01_14:18,2018/06/01_16:10
PaoloG76,MID12,2018/07/01_14:18,2018/07/01_16:10
User1,MID12,2018/07/01_14:18,2018/07/01_16:10
User2,MID12,2009/08/01_14:18,2009/08/01_16:10
PaoloG76,MID10,2018/06/01_14:18,2018/06/01_16:10
PaoloG76,MID10,2018/07/01_14:18,2018/07/01_16:10
User1,MID10,2018/07/01_14:18,2018/07/01_16:10
User2,MID10,2018/08/01_14:18,2018/08/01_16:10
User2,MID124,2001/08/01_14:18,2018/08/01_16:10
User2,MID124,2019/08/01_14:18,2019/08/01_16:10
User2,MID124,2019/08/01_14:18,2019/08/01_16:10
User1,MID12,2019/08/01_14:18,2019/08/01_16:10
User2,MID12,2019/08/01_14:18,2019/08/01_16:10

```

Figure 17 Fichiers de travail

2- Code source Hadoop MapReduce

The image shows two screenshots of a terminal window. The top screenshot shows the execution of a script to start Hadoop daemons. The bottom screenshot shows the compilation of the MapReduce code.

```

omar@omar-HP-Pavillon: ~
omar@omar-HP-Pavillon:~$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as omar in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [omar-HP-Pavillon]
Starting resourcemanager
Starting nodemanagers
omar@omar-HP-Pavillon:~$ jps
13712 SecondaryNameNode
13284 NameNode
13927 ResourceManager
14504 Jps
12235 org.eclipse.equinox.launcher_1.6.500.v20230717-2134.jar
13452 DataNode
14077 NodeManager
omar@omar-HP-Pavillon:~$

omar@omar-HP-Pavillon: ~/eclipse-workspace/Hadoop_MapReduce/src
omar@omar-HP-Pavillon:~/eclipse-workspace/Hadoop_MapReduce/src$ javac -classpath $(hadoop classpath) ourprojct.java
omar@omar-HP-Pavillon:~/eclipse-workspace/Hadoop_MapReduce/src$ ls
Movies.txt  'ourprojctSourprojctMapper.class'  'ourprojctSourprojctReducer.class'  ourprojct.class  ourprojct.java  Users.txt  WatchedMovies.txt
omar@omar-HP-Pavillon:~/eclipse-workspace/Hadoop_MapReduce/src$ jar cf ourprojct.jar 'ourprojctSourprojctMapper.class' 'ourprojctSourprojctReducer.class' ourprojct.class
omar@omar-HP-Pavillon:~/eclipse-workspace/Hadoop_MapReduce/src$ ls
Movies.txt  'ourprojctSourprojctMapper.class'  'ourprojctSourprojctReducer.class'  ourprojct.class  ourprojct.jar  ourprojct.java  Users.txt  WatchedMovies.txt

```

```
omar@omar-HP-Pavilion: ~/eclipse-workspace/Hadoop_MapReduce/src
omar@omar-HP-Pavilion: ~/eclipse-workspace/Hadoop_MapReduce/src$ hadoop jar ourproj.jar ourproj /in/* /outpart1
2023-12-10 18:43:40,216 INFO client.DefaultNoHARMAutoProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2023-12-10 18:43:40,587 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2023-12-10 18:43:40,635 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/omar/.staging/job_1702227518177_0001
2023-12-10 18:43:41,502 INFO input.FileInputFormat: Total input files to process : 6
2023-12-10 18:43:41,632 INFO mapreduce.JobSubmitter: number of splits:6
2023-12-10 18:43:41,864 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1702227518177_0001
2023-12-10 18:43:41,864 INFO mapreduce.JobSubmitter: Executing with tokens: []
2023-12-10 18:43:42,057 INFO conf.Configuration: resource-types.xml not found
2023-12-10 18:43:42,058 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2023-12-10 18:43:42,484 INFO impl.YarnClientImpl: Submitted application application_1702227518177_0001
2023-12-10 18:43:42,529 INFO mapreduce.Job: The url to track the job: http://omar-HP-Pavilion:8088/proxy/application_1702227518177_0001/
2023-12-10 18:43:42,530 INFO mapreduce.Job: Running job: job_1702227518177_0001
2023-12-10 18:43:49,630 INFO mapreduce.Job: Job job_1702227518177_0001 running in uber mode : false
2023-12-10 18:43:49,631 INFO mapreduce.Job: map 0% reduce 0%
2023-12-10 18:44:00,800 INFO mapreduce.Job: map 100% reduce 0%
2023-12-10 18:44:06,843 INFO mapreduce.Job: map 100% reduce 100%
2023-12-10 18:44:07,857 INFO mapreduce.Job: Job job_1702227518177_0001 completed successfully
2023-12-10 18:44:07,996 INFO mapreduce.Job: Counters: 54
File System Counters
  FILE: Number of bytes read=78
  FILE: Number of bytes written=1935190
```

Figure 18 Code source Hadoop MapReduce

=> Mapper class

```
// Mapper Class
public static class OurProjMapper extends Mapper<Object, Text, Text, Text> {
    private Text movieId = new Text();
    private Text username = new Text();

    // Map method processes each input record
    public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
        // Split the CSV record into tokens
        String[] tokens = value.toString().split(",");

        // Assuming the format is "Username, MID, StartTimestamp, EndTimestamp"
        if (tokens.length == 4) {
            String startTimestamp = tokens[2].trim();
            // Check if the start timestamp is in the year 2019
            if (isInYear2019(startTimestamp)) {
                movieId.set(tokens[1].trim());
                username.set(tokens[0].trim());
                // Emit movieId as key and username as value
                context.write(movieId, username);
            }
        }
    }
}
```

Figure 19 Mapper Class

=> Reduce class

```
// Reducer Class
public static class OurProjetReducer extends Reducer<Text, Text, Text, Text> {
    private Text result = new Text();

    // Reduce method processes the grouped key-value pairs
    public void reduce(Text key, Iterable<Text> values, Context context)
        throws IOException, InterruptedException {
        Set<String> uniqueUsers = new HashSet<>();

        // Iterate through values and add unique usernames to the set
        for (Text value : values) {
            uniqueUsers.add(value.toString());
        }

        // If there is only one unique user, emit the key-value pair
        if (uniqueUsers.size() == 1) {
            result.set(uniqueUsers.iterator().next());
            context.write(key, result);
        }
    }
}
```

Figure 20 Reduce class

Résultat:

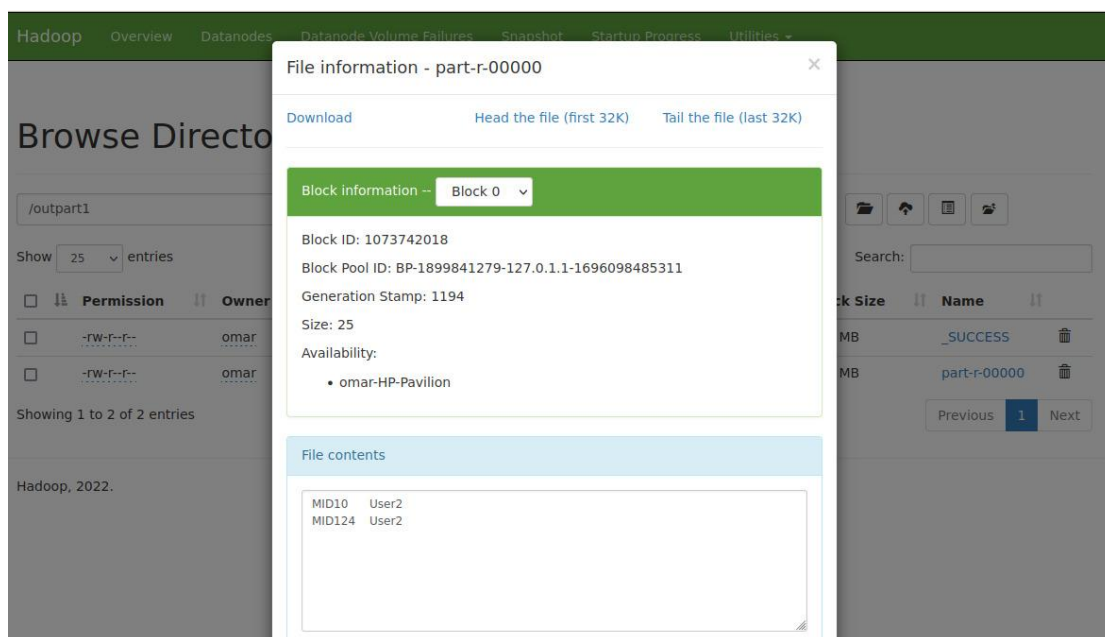


Figure 21 Resultat partie 1

3- Code source Spark et RDD

Question 1:

Films qui ont été regardés fréquemment mais seulement pendant un an au cours des cinq dernières années.


```

1 from pyspark import SparkContext
2 from datetime import datetime
3
4 # Initialisation du contexte Spark
5 sc = SparkContext(appName="Films_Frequents_Cinq_Dernieres_Annees")
6
7 # Chargement du fichier WatchedMovies.txt en tant que RDD
8 watched_movies_rdd = sc.textFile("WatchedMovies.txt")
9
10 # Filtrage des lignes liées aux cinq dernières années
11 filtered_movies_rdd = watched_movies_rdd.filter(lambda line: "2015/09/17" <= line.split(",")[2] <= "2020/09/16")
12
13 # Transformation en paire clé-valeur avec (MID, (Année, 1)) pour chaque ligne
14 movies_year_count_rdd = filtered_movies_rdd.map(lambda line: ((line.split(",")[1], line.split(",")[2][:4]), 1))
15
16 # Réduction par clé pour compter le nombre de fois que chaque film a été regardé chaque année
17 movies_year_count_sum_rdd = movies_year_count_rdd.reduceByKey(lambda x, y: x + y)
18
19 # Filtrage des films qui ont été regardés au moins 1000 fois au cours d'une année
20 selected_movies_rdd = movies_year_count_sum_rdd.filter(lambda x: x[1] >= 3)
21
22 # Transformation pour obtenir le format (MID, Année)
23 result_rdd = selected_movies_rdd.map(lambda x: (x[0][0], x[0][1]))
24
25 # Affichage des résultats
26 result_rdd.foreach(print)
27
28 # Enregistrement des résultats dans un fichier de sortie (vous pouvez ajuster le chemin selon vos besoins)
29 result_rdd.saveAsTextFile("out_part2_Q2")
30
31 # Arrêt du contexte Spark
32 sc.stop()
33

```

Figure 22 Code question1 spark

```

code_Q1.py  part-00001
out_part2_Q2 > part-00000
1 ('MID10', '2018')
2

```

Figure 23 Output question1 spark

Question2:

Film le plus populaire depuis au moins deux ans.

```

1 from pyspark import SparkContext, SparkConf
2
3 # Create a Spark configuration and set the application name
4 conf = SparkConf().setAppName("MostPopularMoviesByYear")
5 sc = SparkContext(conf=conf)
6
7 # Load the WatchedMovies data from the input file
8 watched_movies = sc.textFile("WatchedMovies.txt")
9
10 # Split each line into fields
11 # Assuming the format is "Username, MID, StartTimestamp, EndTimestamp"
12 # Adjust the split logic based on your actual data format
13 watched_movies_data = watched_movies.map(lambda line: line.split(","))
14
15 # Map the data to (MID, (Year, User)) pairs
16 movie_year_user_pairs = watched_movies_data.map(lambda fields: ((fields[1], fields[2].split("/")[-1]), fields[0]))
17
18 # Remove duplicate users in a specific year for each movie
19 distinct_users_per_year = movie_year_user_pairs.distinct()
20
21 # Map to (MID, 1) to count occurrences of each movie in each year
22 movie_count = distinct_users_per_year.map(lambda x: ((x[0][0], x[0][1]), 1))
23

```

```

21 # Map to (MID, 1) to count occurrences of each movie in each year
22 movie_count = distinct_users_per_year.map(lambda x: ((x[0][0], x[0][1]), 1))
23
24 # Reduce by key to get the count of distinct users for each movie in each year
25 movie_user_count = movie_count.reduceByKey(lambda x, y: x + y)
26
27 # Map to (Year, (MID, Count)) to find the most popular movie for each year
28 most_popular_movies = movie_user_count.map(lambda x: (x[0][1], (x[0][0], x[1])))
29
30 # Reduce by key to get the final result, selecting the most popular movie for each year
31 result = most_popular_movies.reduceByKey(lambda x, y: x if x[1] > y[1] else y).map(lambda x: (x[0], x[1][0]))
32
33 # Collect the result to the driver program
34 collected_result = result.collect()
35
36 # Count the occurrences of each MID
37 mid_counts = {}
38 for _, mid in collected_result:
39     mid_counts[mid] = mid_counts.get(mid, 0) + 1
40
41 # Filter MIDs that are present in at least 2 years
42 filtered_mids = [mid for mid, count in mid_counts.items() if count >= 2]
43
44 # Filter the collected result based on the selected MIDs
45 final_result = [(year, mid) for year, mid in collected_result if mid in filtered_mids]
46
47 # Save the result to the output directory
48 sc.parallelize(final_result).map(lambda x: f"{x[1]}").saveAsTextFile("OutPart2_Q2")
49
50 # Stop the Spark context
51 sc.stop()

```

Figure 24 Code question2 spark

part-00005	×
out_part2_Q2	>
1	MID124
2	

part-00002	×
out_part2_Q2	>
1	MID12
2	

Figure 25 Output question2 spark

FIN