

# Recherche & Développement autour de la problématique du Big Data



VIGNERON

Joris



# ☐ Sommaire

1. Introduction
2. Présentation de l'entreprise
3. Sujet du stage
4. Etat de l'art
5. Découverte et expérimentation
6. Application
7. Conclusion

# Introduction

- Choix de l'Observatoire astronomique
- 2 raisons ayant motivé ce choix

# Présentation de l'entreprise

- Stage réalisé à l'Observatoire astronomique de Strasbourg
- Construit en 1881, possède la 3<sup>ème</sup> plus grande lunette de France par sa taille
- Il se doit de contribuer aux progrès de la connaissance



Observatoire astronomique  
de Strasbourg

# Présentation de l'entreprise – Les équipes de recherche

- Galaxies : réalise des recherches sur la formation des galaxies
- Hautes énergies : s'intéresse aux sources émettrice de rayon X
- Le CDS : développe des services d'accès aux données



# Présentation de l'entreprise - Les services du CDS



Base de données de référence pour la bibliographie des objets astronomiques

Plus de 40 ans de données, informations sur plus de 8 Millions d'objets

Base de données qui regroupe plus de 13 000 catalogues d'objets célestes. 300 000 requêtes par jour en moyenne.

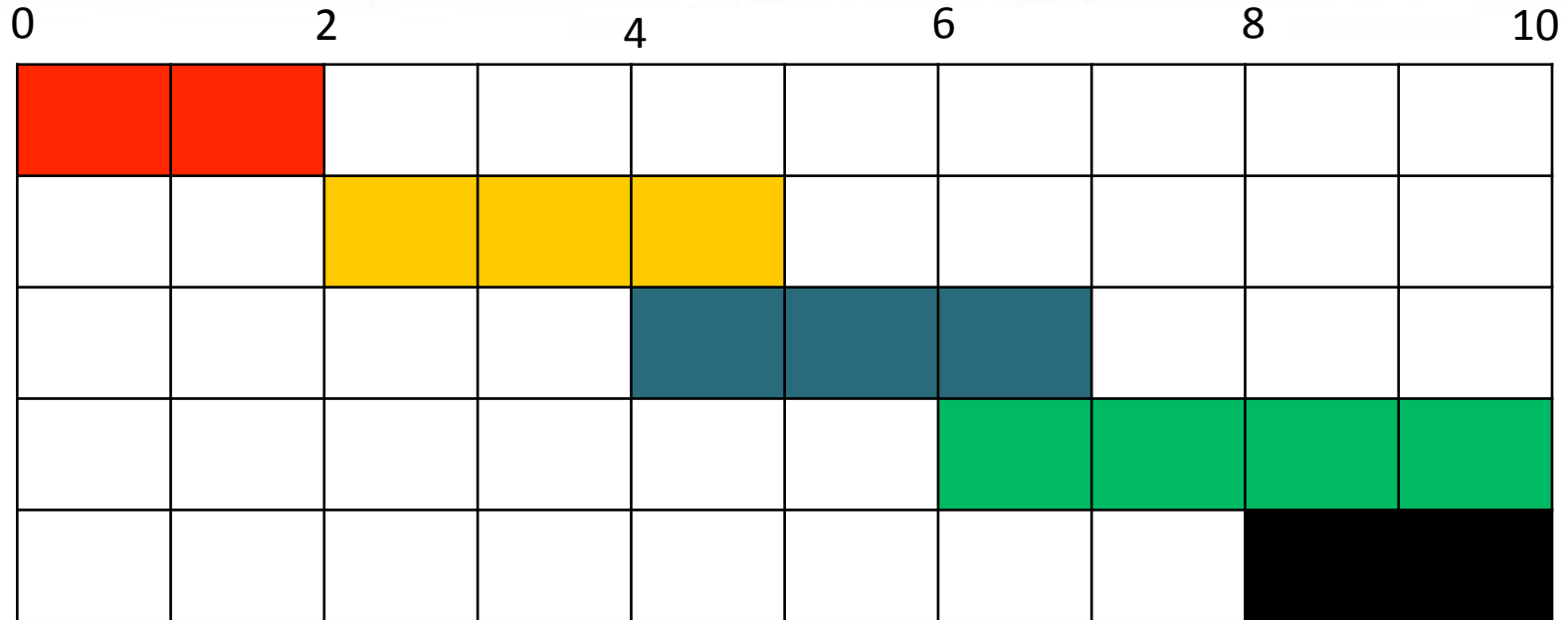


Atlas interactif du ciel permettant de visualiser des images astronomiques

# Sujet du stage

- Problématique du Big Data et notamment les technologies qui gravitent autour
- Enjeux du stage :
  - Réalisation de diverses expérimentations dans le but d'éprouver les outils à disposition
  - Mise en œuvre d'un cas d'utilisation de l'Observatoire

# Sujet du stage - Déroulement



Etat de l'art



Essais avec Spark



Rédaction du rapport



Essais avec Hadoop



Réalisation de l'application



# Etat de l'art

- Début du stage : découvrir les technologies du Big Data
- Apprentissage des différentes facettes de cette technologie



# Etat de l'art - MapReduce

- Patron d'architecture de programmation
- Manipulation de grandes quantités de données
- Fonctionnement en 2 phases : phase Map et phase Reduce

# Etat de l'art - Hadoop

- Framework Java Open Source développé par la fondation Apache



- Hadoop Distributed File System, HDFS est le système de fichier distribué d'Hadoop
- Basé sur le patron MapReduce
- Différentes distributions implémentant ce framework

# Etat de l'art - Spark

- Framework développé par AMPLab de l'université de Berkeley
- Modèle de programmation plus simple et temps d'exécution 100 fois plus courts
- Le RDD, « Resilient Distributed Dataset », est une abstraction de collection sur laquelle est réalisée des opérations de manières distribuées
- Programmation différente d'Hadoop



# Etat de l'art – Pig & Hive

- Infrastructure construite par-dessus Hadoop
- Utilise PigLatin, un langage de script
- Hive, permet de réaliser des requêtes en HiveQL, un langage proche du SQL
- Très bon outil pour développer de rapide job MapReduce



# Découverte et expérimentation

- Objectif prendre en main le framework Hadoop
- Choix de la distribution HortonWorks
- Installation d'une distribution en local avec Eclipse



## Hortonworks Sandbox with HDP 2.2

[Leave Feedback](#)

Component	Version
Map	2.6.1.0
HDP	2.2.4
Hadoop	2.6.0
Pig	0.14.0
Hive-Metastore	0.14.0
Cloudera	4.1.0
Ambari	2.5.10 <a href="#">Details</a>
Yarn	2.6.0
Tez	0.8.4
Storm	0.9.3
Falcon	0.8.0
Sandbox Build	114638 (0-25 03-04-15)





# Découverte et expérimentation -Hive

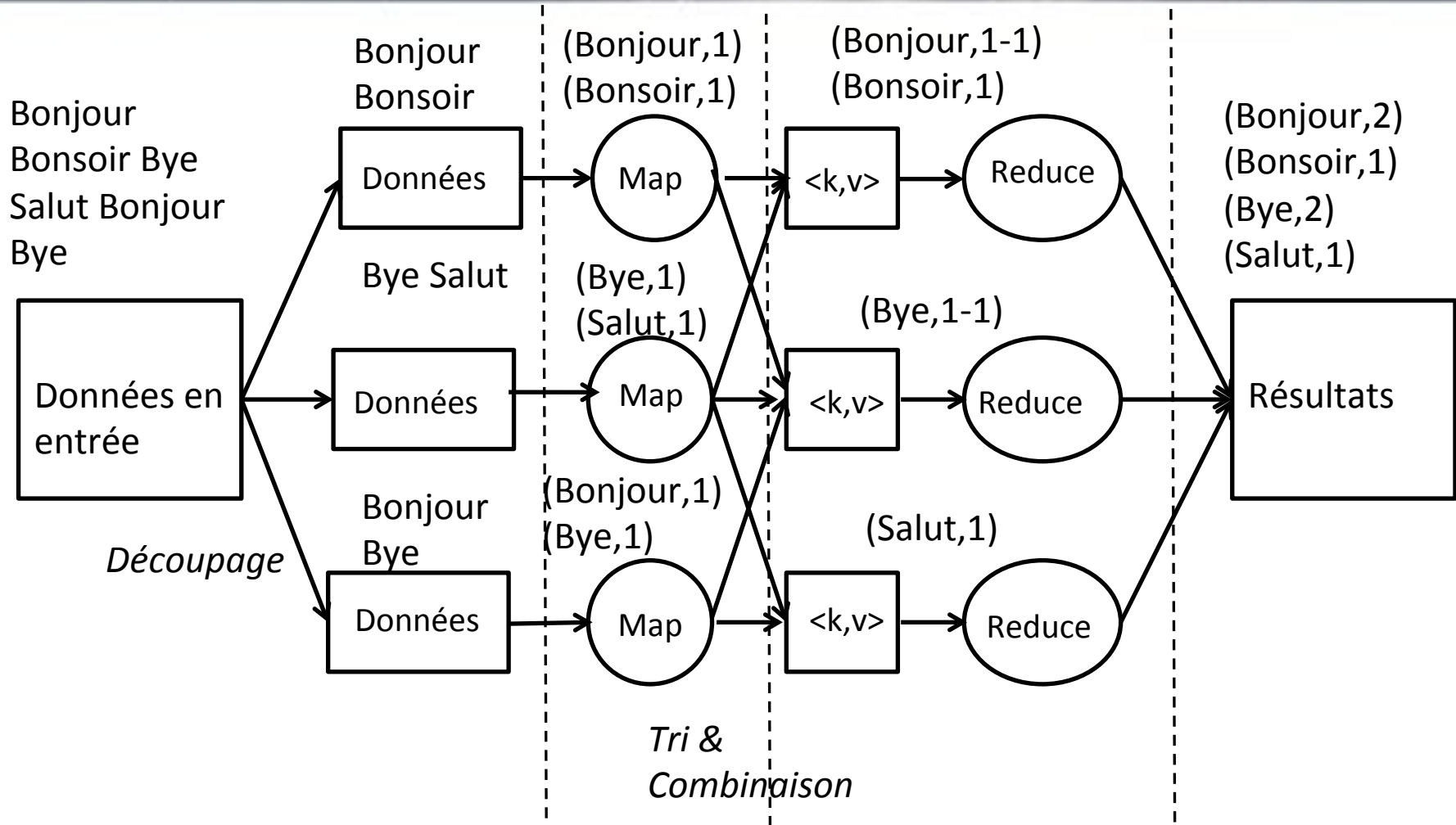
- Simplification d'un job MapReduce
- Utilisation d'un langage HiveQL très proche du SQL

# Bilan sur Hive et Pig

- Accessible pour des programmeurs non formés à Hadoop
- Inconvénients :
  - Presque inefficace pour des programmes trop complexes
  - Ne supporte pas tous les standards attendus

# Découverte et expérimentation – Essais avec Hadoop

- Utilisation de Maven pour créer les projets
- Programmeur doit coder : une fonction *map()*, une fonction *reduce()*, et un driver
- Les données chargées sont découpées sous forme de bloc puis distribuées sur le cluster



# Découverte et expérimentation – Essais avec Spark

- Chargement des données dans des javaRDD
- Utilisation de méthode emprunter au langage fonctionnel pour réaliser des traitements sur les JavaRDD
- Utilisation des classes Function, Function2,... permettant de définir les types d'entrées et de sorties

```

class PairRDD[String, String] extends RDD[String] with Serializable {
  private val f: (String, String) => String = (a, b) => a + b

  public Tuple2[String, String] call(String a) {
    // Je décompose la ligne pour récupérer chaque des valeurs
    String[] t = a.split(",")
    // Je stocke la valeur ra et de et je les converti en double
    Double ra = Double.parseDouble(t[0])
    Double de = Double.parseDouble(t[1])
    // Je calcule la valeur de l'itération
    Long hald = nextId(ra, de)
    return new Tuple2[String, String](hald, callString(t, a))
  }
}

```

Application d'une fonction, dont le paramètre est une classe dont on peut définir les types d'entrées et de sorties, puis application sur chaque élément de notre JavaRDD d'un traitement correspondant à la fonction call contenu dans la classe interne.

# Application

- Détermination d'une application à réaliser pour l'Observatoire
- Réalisation d'une application de CrossMatch



# Application - CrossMatch

- Réalisation d'une jointure entre 2 catalogues
- Calcul d'un identifiant servant de clé de jointure
- Filtration des sources en calculant la distance entre elles

# Application - Description des sources

- 2 Catalogues :
  - Tycho : 2 539 913 valeurs
  - HipMain : 119 218 valeurs

```
45.084243,+0.260814,10.723,-1  
45.084843,+0.260803,10.977,-1  
45.104016,+0.300000,10.500,-1  
45.140714,+0.300004,10.304,-1  
45.182794,+0.300000,10.100,-1  
45.112794,+0.300014,10.204,-1  
45.010100,+0.301004,10.100,-1  
45.004000,+0.470000,9.970,-1  
45.074004,+0.470007,10.300,-1
```

```
3, 9.27, 0.00379707, -19.48882768, 21.9  
3, 8.41, 0.00000700, 28.80018000, 3.82  
4, 9.04, 0.00000007, -41.89914432, 7.76  
5, 9.56, 0.00000004, -49.85120448, 2.87  
6, 12.01, 0.01014144, 3.74048830, 28.8  
7, 9.66, 0.00204001, 10.00000000, 27.76  
8, 9.04, 0.02000000, 28.80018000, 8.17  
9, 9.04, 0.00000000, 49.80000000, 4.80  
10, 8.00, 0.00000000, -50.00000000, 10.74
```

# Application – Réalisation avec Hadoop

- Réalisation de recherches pour trouver une façon de faire une jointure
- Construction du code :
  - 2 Mappers, un par catalogue
  - 1 Reducer
  - 1 Driver

# Application - HipMainJoinMapper

- Découpage de la ligne courante pour obtenir valeurs de ra et de
- Calcul de l'identifiant à partir de ces 2 valeurs
- Création du couple <clé,valeur> :
  - Clé -> l'identifiant
  - valeur -> la ligne courante + une lettre

```
@Override
public void map(Object key, Object value, OutputCollector context,
    Reporter arg) throws IOException {
    String[] parsed = value.toString().split(",");
    if(parsed.length==4){
        String ra = parsed[2];
        String dec = parsed[3];
        double raDeg = Double.parseDouble(ra);
        double decDeg = Double.parseDouble(dec);
        //calcul de l'id
        long id = nextId(raDeg,decDeg);
        //on passe l'id comme clé du couple
        outkey.set(id+"");
        //on passe la valeur de la ligne -> à pour désigner la src d'origine
        outvalue.set("R" + value.toString());
        //on passe le couple valeur en sortie
        context.collect(outkey,outvalue);

        // traite les pixels voisins
    }
    try {
        for (long nei : nb.neighbours(4)) {
            outkey.set(nei+"");
            context.collect(outkey,outvalue);
        }
    } catch (Exception e) {
        // TODO Auto-generated catch block
        e.printStackTrace();
    }
}
```

# Application - CrossMatchReducer

- En entrées : un couple <clé,Collection<valeur>>
- Ces valeurs viennent des 2 catalogues, identification par la lettre
- Parcours de la Collection<valeur> pour ranger chaque valeur dans la liste correspondante
- Réalisation de la jointure en parcourant les 2 listes, filtration des données avant écriture

```
private void executeJoinLogic(OutputCollector context) throws IOException {
    if(!listA.isEmpty() && !listB.isEmpty()){
        for(Text a : listA){
            // EXTRACT POS A
            String[] a = a.toString().split(",");

            double val = Double.parseDouble(a[0]);
            double dec1 = Double.parseDouble(a[1]);

            for(Text b : listB) {
                // EXTRACT POS B
                String[] b = b.toString().split(",");
                double val2 = Double.parseDouble(b[0]);
                double dec2 = Double.parseDouble(b[1]);
                double distance = haversine(val,dec1,val2,dec2);
                if (distance < matchDistance) {
                    context.collect(a,b);
                }
            }
        }
    }
}

@Override
public void reduce(Object key, Iterator values, OutputCollector context,
    Reporter arg) throws IOException {
    listA.clear();
    listB.clear();
    while(values.hasNext()){
        Text tmp = (Text) values.next();
        if(tmp.charAt(0) == 'A'){
            listA.add(new Text(tmp.toString().substring(1)));
        }else if(tmp.charAt(0) == 'B'){
            listB.add(new Text(tmp.toString().substring(1)));
        }
    }
    try {
        executeJoinLogic(context);
    } catch (InterruptedException e) {
        e.printStackTrace();
    }
}
```

# Application - Résultat

- En entrée on a donc :
  - Tycho : 2,5 Millions valeurs
  - HipMain : 120 000 valeurs
- En sortie : 119 167 valeurs

# Conclusion

- Une expérience enrichissante
- Difficultés rencontrées
- Conclusion sur Hadoop et Spark



Merci de m'avoir écouté, avez-vous  
des questions ?