



Projet Big Data 2019: GDELT

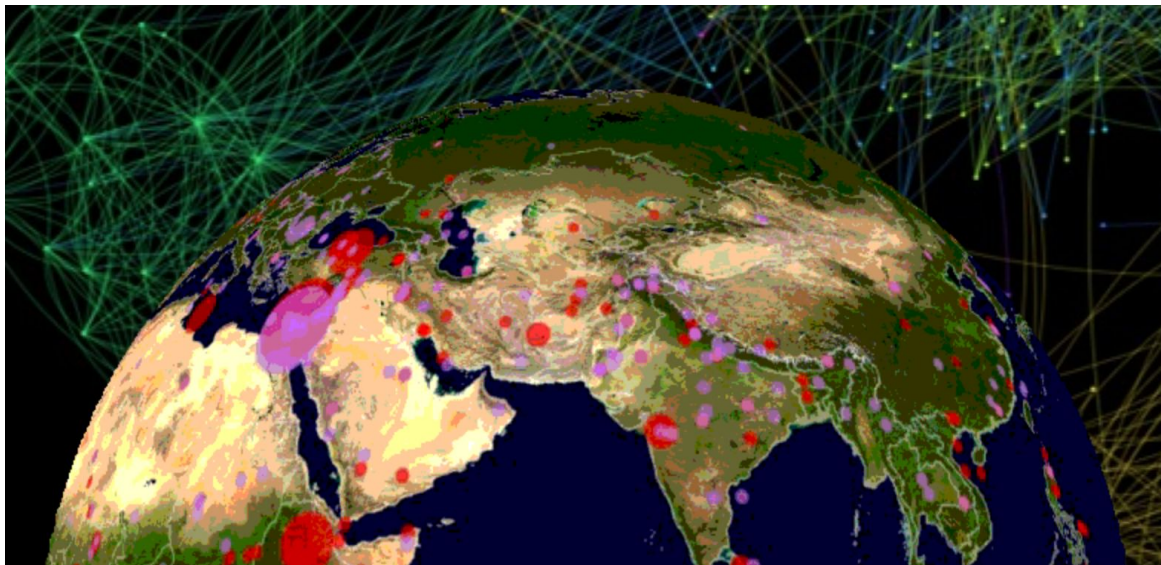
Emmanuel, Hicham, Paul, Jean-Galbert, Ialifinaritra

Sommaire

- Contexte
- Description des données
- Technologies choisies
- Architecture
- Requêtage
- Demo
- Améliorations

Objectif du projet

- **But :** Proposer un système de stockage distribué, résilient et performant sur AWS.
- **Base de données :** <http://data.gdeltproject.org/gdeltv2/>
 - Données d'articles de presse enregistrées toutes les 15 minutes et stockées en 3 CSV.



Description des données

Événements

Evènement_Global_ID

Date

Pays

Nombre d'articles

Mentions

Ton moyen

Mentions

Evènement_Global_ID

Date

Source

Langue

GKG

Nom de la source

Date

Thème

Personne

Locations

Ton

Nombre d'articles

Technologies

- AWS :



amazon
EMR



Amazon
EC2



- Calcul distribué :



- Résilience du stockage des données:

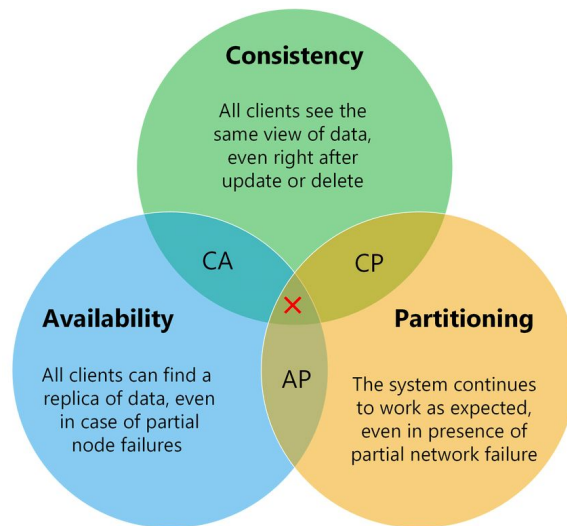


cassandra

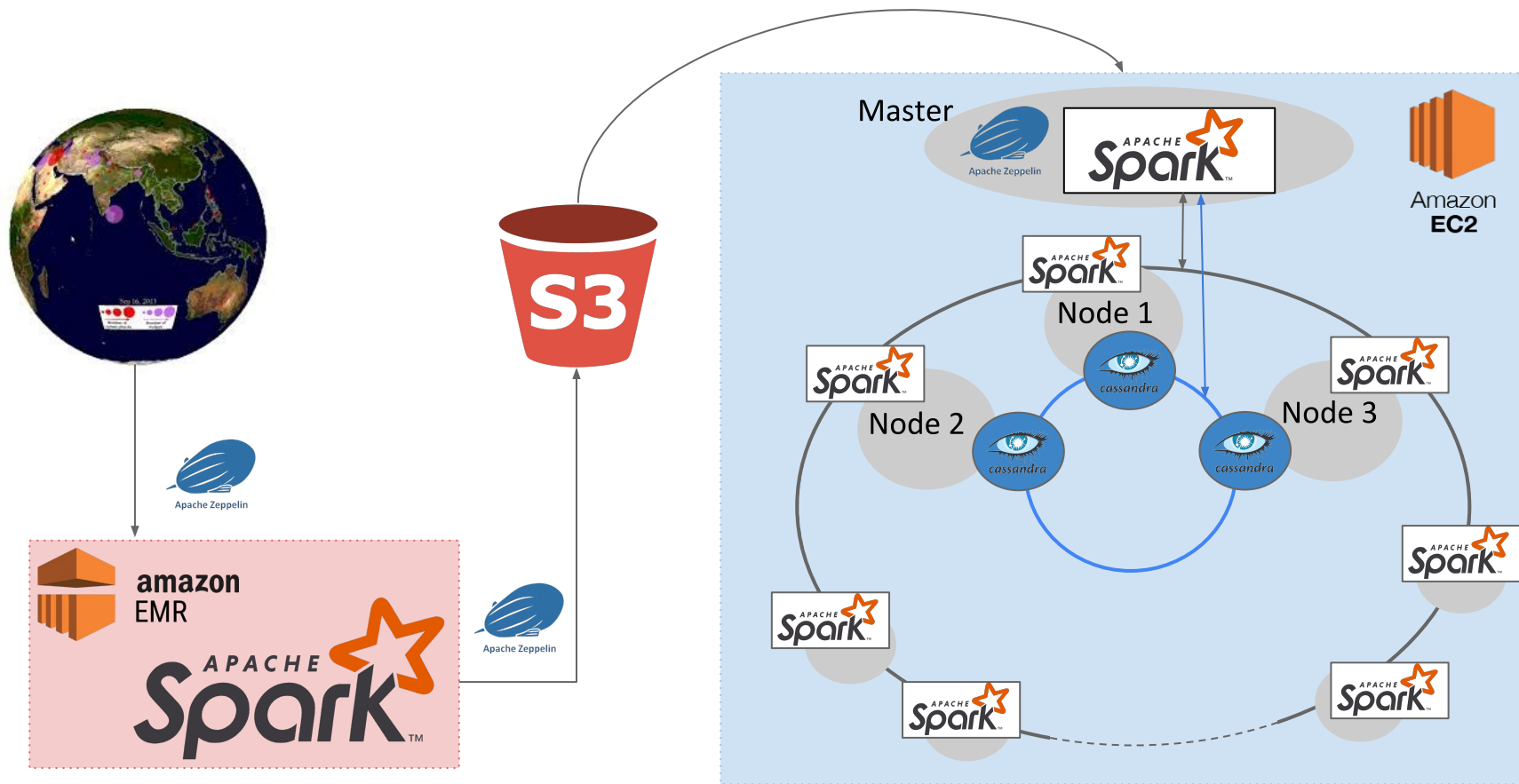
- Notebook:



Apache Zeppelin



Visualisation de l'architecture



Choix des technologies :

Avantages :

- Flexibilité des installations sur EC2
- Configuration minimale de Cassandra
- Documentation fournie sur Cassandra
- Bon outil de visualisation sur Zeppelin

Contraintes :

- Compte Educate !
- Comportement global du système difficile à débbugger

Modélisation

Requête 1 : Afficher le nombre d'articles/événements qu'il y a eu pour chaque triplet (jour, pays de l'évènement, langue de l'article).

Evénements



Mentions

TABLE 1
<u>Jour</u>
<u>Pays</u>
<u>Langue</u>
Mentions
Evénements

Modélisation

Requête 2 : Pour un pays donné en paramètre, afficher les événements qui y ont eu place triés par le nombre de mentions (tri décroissant)

Événements



TABLE 2
<u>Evènement Global ID</u>
<u>Pays</u>
<u>Jour</u>
<u>Mois</u>
<u>Année</u>
Nombre Mentions

Modélisation

Requête 3 : Pour une source de données passée en paramètre ("gkg.SourceCommonName"), afficher les thèmes, personnes, lieux dont les articles de cette source parlent ainsi que le nombre d'articles et le ton moyen des articles (pour chaque thème/personne/lieu)

TABLE 3_1 : Thèmes

<u>Source Commun Nom</u>
<u>Jour</u>
<u>Mois</u>
<u>Année</u>
Thème
Ton moyen
Nombre Articles

TABLE 3_2 : Personnes

<u>Source Commun Nom</u>
<u>Jour</u>
<u>Mois</u>
<u>Année</u>
Personne
Ton moyen
Nombre Articles

TABLE 3_3 : Locations

<u>Source Commun Nom</u>
<u>Jour</u>
<u>Mois</u>
<u>Année</u>
Location
Ton moyen
Nombre Articles

Modélisation

Requête 4 : Dresser la cartographie des relations entre les pays d'après le ton des articles : pour chaque paire (Pays 1, Pays 2), calculer le nombre d'article, le ton moyen (agrégations sur Année/Mois/Jour, filtrage par pays ou carré de coordonnées)

Events

Mentions



TABLE 4
<u>Pays1</u>
<u>Pays2</u>
<u>Jour</u>
<u>Mois</u>
<u>Année</u>
Tonmoyen
NbArticles

Demo

Points d'amélioration

- Utilisation de EMR pour le chargement des données.
- Plus des noeuds Cassandra pour accélérer l'écriture.
- Requêtes optimisées, compilées en JAR pour accélérer le nettoyage et chargement les données