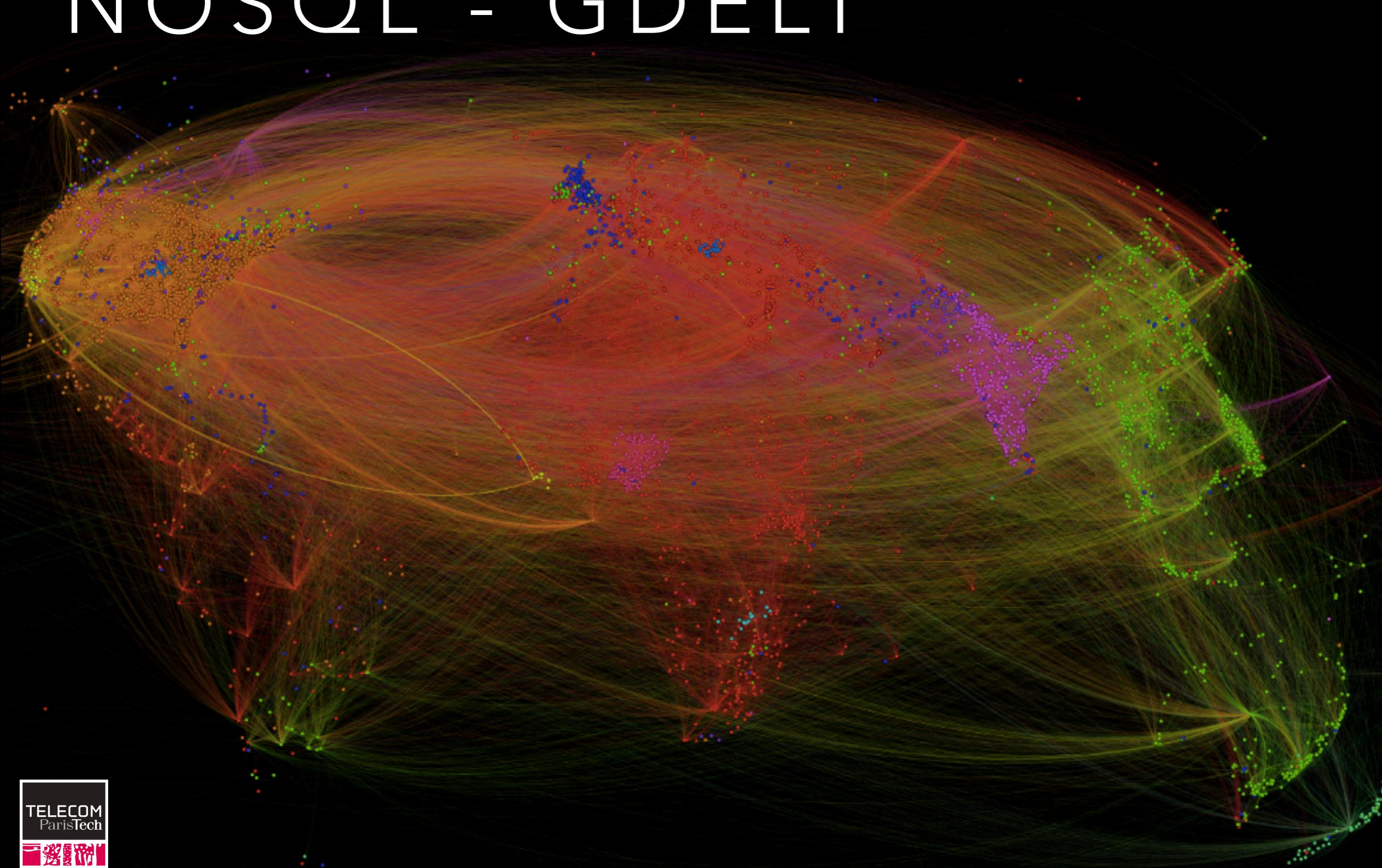


TELECOM PARISTECH - 2019

NOSQL - GDELT



GDELT

SUJET



PROBLÉMATIQUE

- Proposer un système de stockage distribué, résilient et performant sur AWS pour les données de GDELT.

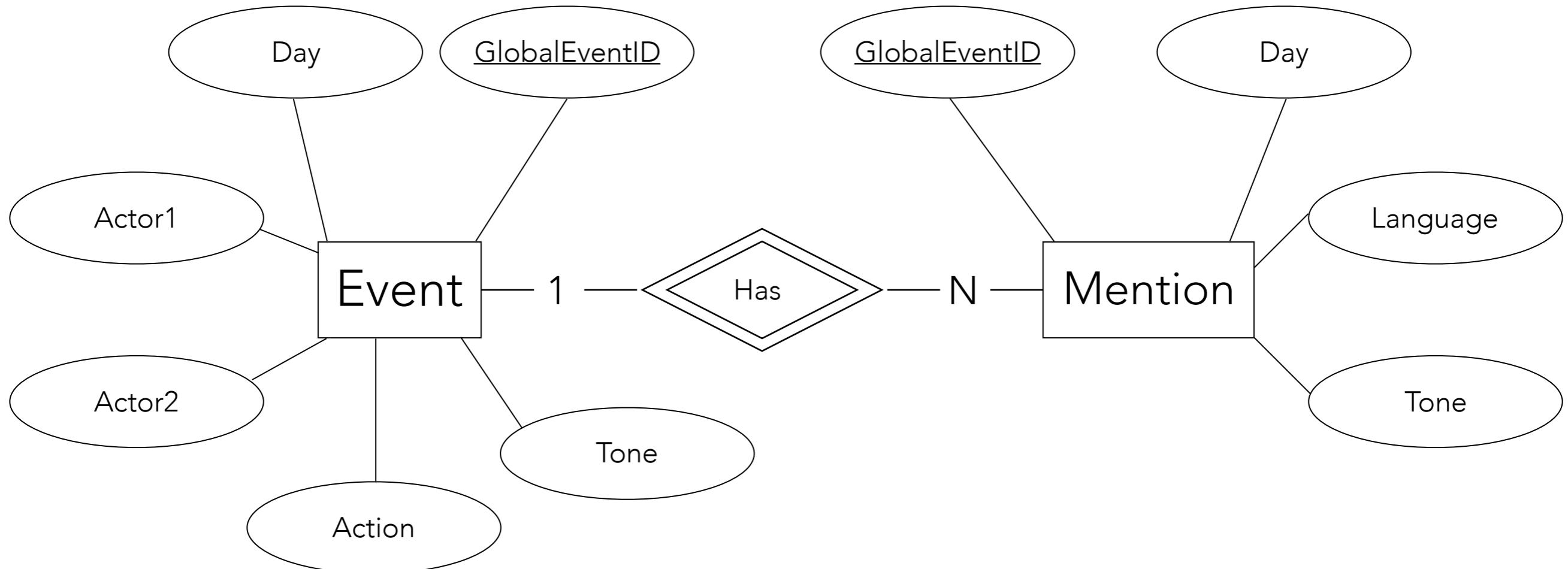
Afficher :

- le nombre d'articles / événements pour chaque (jour, pays de l'événement, langue de l'article)
- pour un acteur(pays/organisation ...) ⇒ afficher les événements qui y font référence
- les sujets (acteurs) qui ont eu le plus d'articles positifs/négatifs (mois, pays, langue de l'article)
- acteurs/pays/organisations qui divisent le plus

LES DONNÉES



MODÈLE CONCEPTUEL

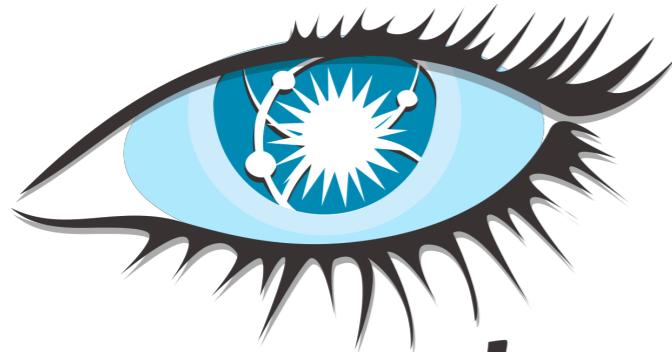


GDELT

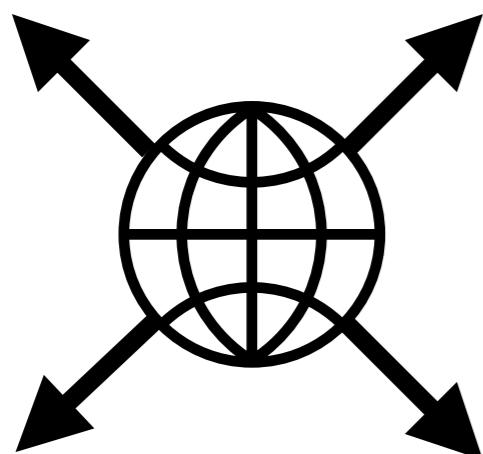
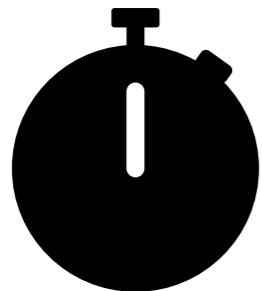
APPROCHE



BASE DE DONNÉES



cassandra

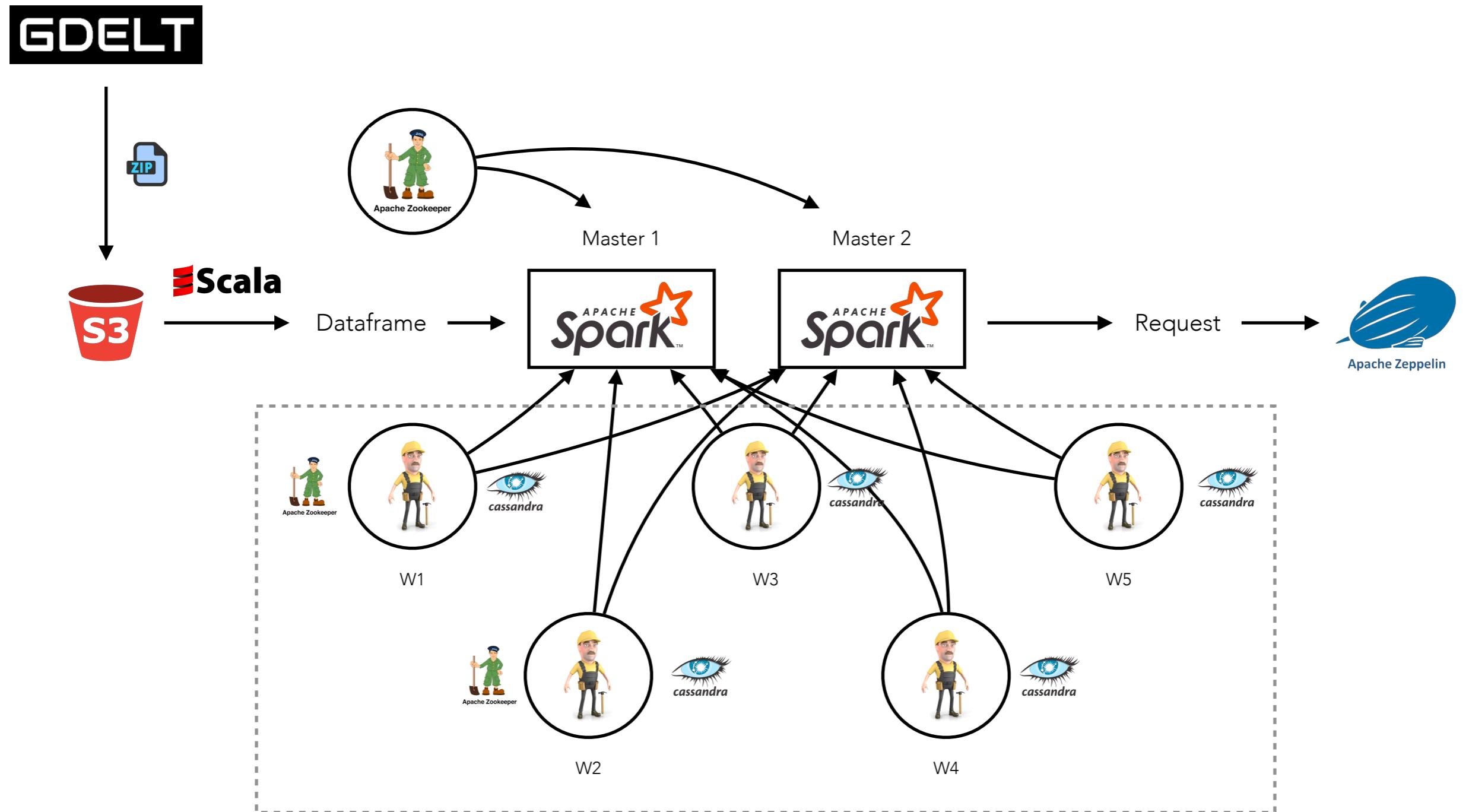


Exemple :

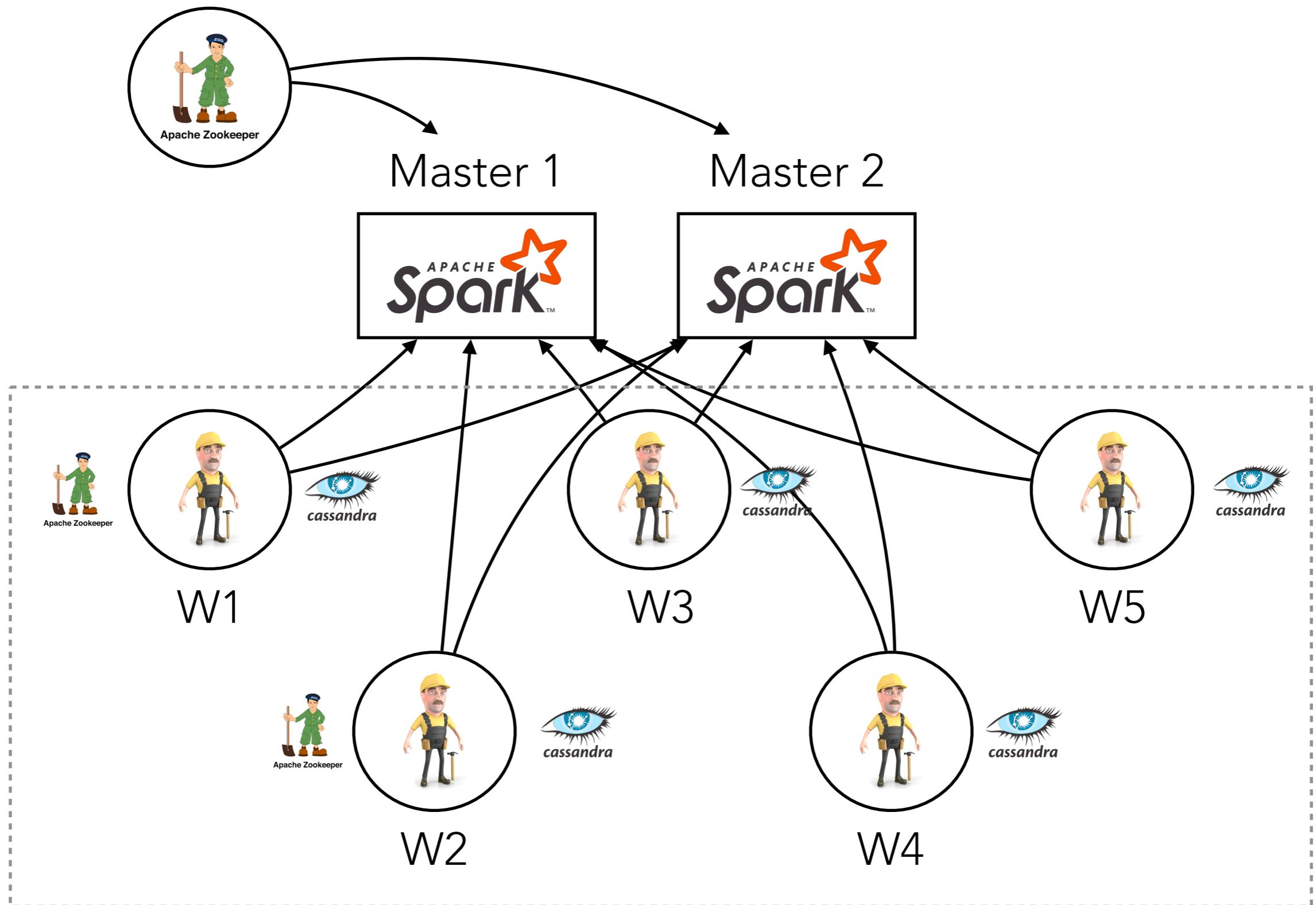
Partition Key	Columns
<u>GlobalEventID</u>	
	Day Language NumArticles
	20180112 English 248'540
	Day Language NumArticles
	20180113 English 129'540
	...

Replication Factor : 3

ARCHITECTURE



ARCHITECTURE



DATA VISUALIZATION

- Zeppelin `z.show(spark.sql(""" SELECT * FROM q3month ORDER BY SumTone DESC LIMIT 100 """))`
- Helium Zeppelin Leaflet
- Scala + Spark

GDELT

DIFFICULTÉS



DIFFICULTÉS

- Problème avec EC2 la veille de la présentation

```
java.lang.NoSuchMethodError: com.amazonaws.services.s3.transfer.TransferManager.<init>(Lcom/amazonaws/services/s3/AmazonS3;Ljava/util/concurrent/ThreadPoolExecutor;)V
  at org.apache.hadoop.fs.s3a.S3AFileSystem.initialize(S3AFileSystem.java:287)
  at org.apache.hadoop.fs.FileSystem.createFileSystem(FileSystem.java:2669)
  at org.apache.hadoop.fs.FileSystem.access$200(FileSystem.java:94)
  at org.apache.hadoop.fs.FileSystem$Cache.getInternal(FileSystem.java:2703)
  at org.apache.hadoop.fs.FileSystem$Cache.get(FileSystem.java:2685)
  at org.apache.hadoop.fs.FileSystem.get(FileSystem.java:373)
  at org.apache.hadoop.fs.Path.getFileSystem(Path.java:295)
  at org.apache.hadoop.mapreduce.lib.input.FileInputFormat.setInputPaths(FileInputFormat.java:500)
  at org.apache.hadoop.mapreduce.lib.input.FileInputFormat.setInputPaths(FileInputFormat.java:469)
  at org.apache.spark.SparkContext$$anonfun$binaryFiles$1.apply(SparkContext.scala:921)
  at org.apache.spark.SparkContext$$anonfun$binaryFiles$1.apply(SparkContext.scala:916)
  at org.apache.spark.rdd.RDDOperationScope$.withScope(RDDOperationScope.scala:151)
  at org.apache.spark.rdd.RDDOperationScope$.withScope(RDDOperationScope.scala:112)
  at org.apache.spark.SparkContext.withScope(SparkContext.scala:693)
  at org.apache.spark.SparkContext.binaryFiles(SparkContext.scala:916)
... 62 elided
```

- Passage sur EMR
- Préparation des données avant passage dans Cassandra à refaire

GDELT

PERFORMANCES



PERFORMANCES

- Requête d'un mois en Spark SQL :

The screenshot shows a Jupyter Notebook cell with the following content:

```
z.show(spark.sql(""" SELECT * FROM q3 ORDER BY SumTone ASC LIMIT 100 """))
```

Below the code, there is a visualization toolbar with icons for different chart types (grid, bar, pie, map, line, scatter) and a download button. To the right of the toolbar is the text "settings ▾".

The main area displays a table with the following data:

Language	ActionCountry	Month	Actor1Country	Actor1Code	SumTone
eng	US	201801	US	COP	-1365577.0567040334
eng	US	201801	US	JUD	-1037184.958702697
eng	US	201801	US	GOV	-1029943.3650900614
eng		201801		COP	-696337.9742912925
ara	IS	201801	IS	ISR	-694516.3164537457
eng	IS	201801	IS	ISR	-599651.2468628696
eng	IR	201801	IR	IRN	-476389.1384875841
eng	RS	201801	RS	RUS	-460193.67098307423

Took 10 min 48 sec. Last updated by anonymous at January 23 2019, 10:41:24 AM.

- Requête d'un mois en Cassandra :

Took 10 min 48 sec. Last updated by anonymous at January 23 2019, 10:41:24 AM.

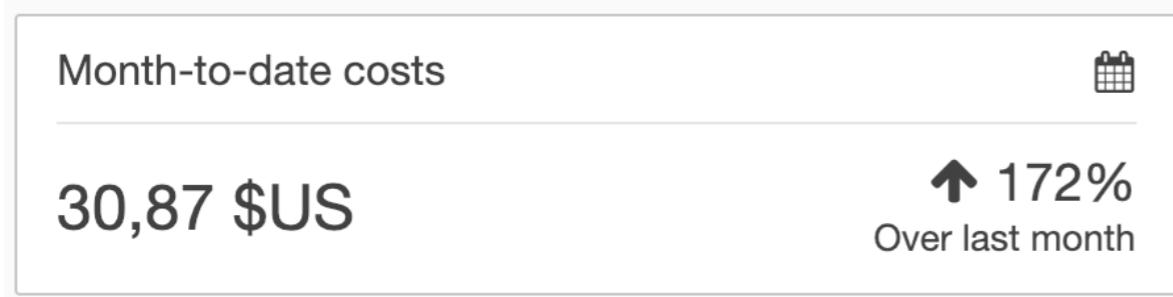
GDELT

BUDGET

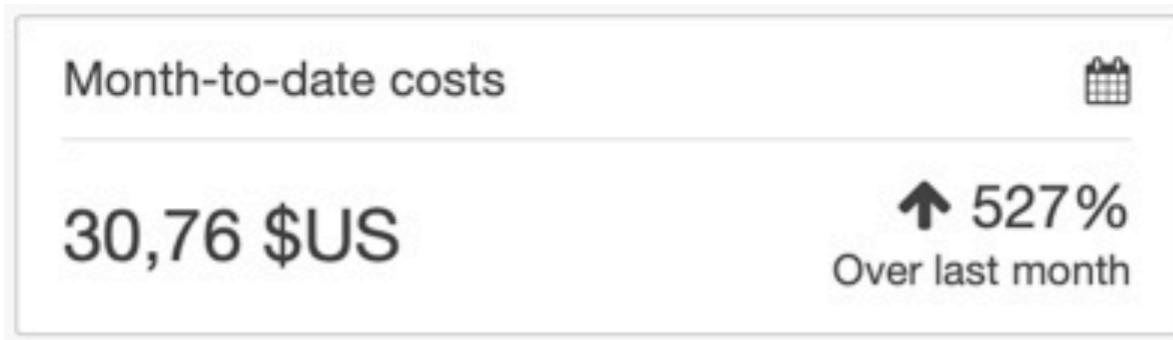


BUDGET

- Charger et stocker les données dans S3, pre-processing vers Cassandra :



- Architecture EC2 avant shift EMR :



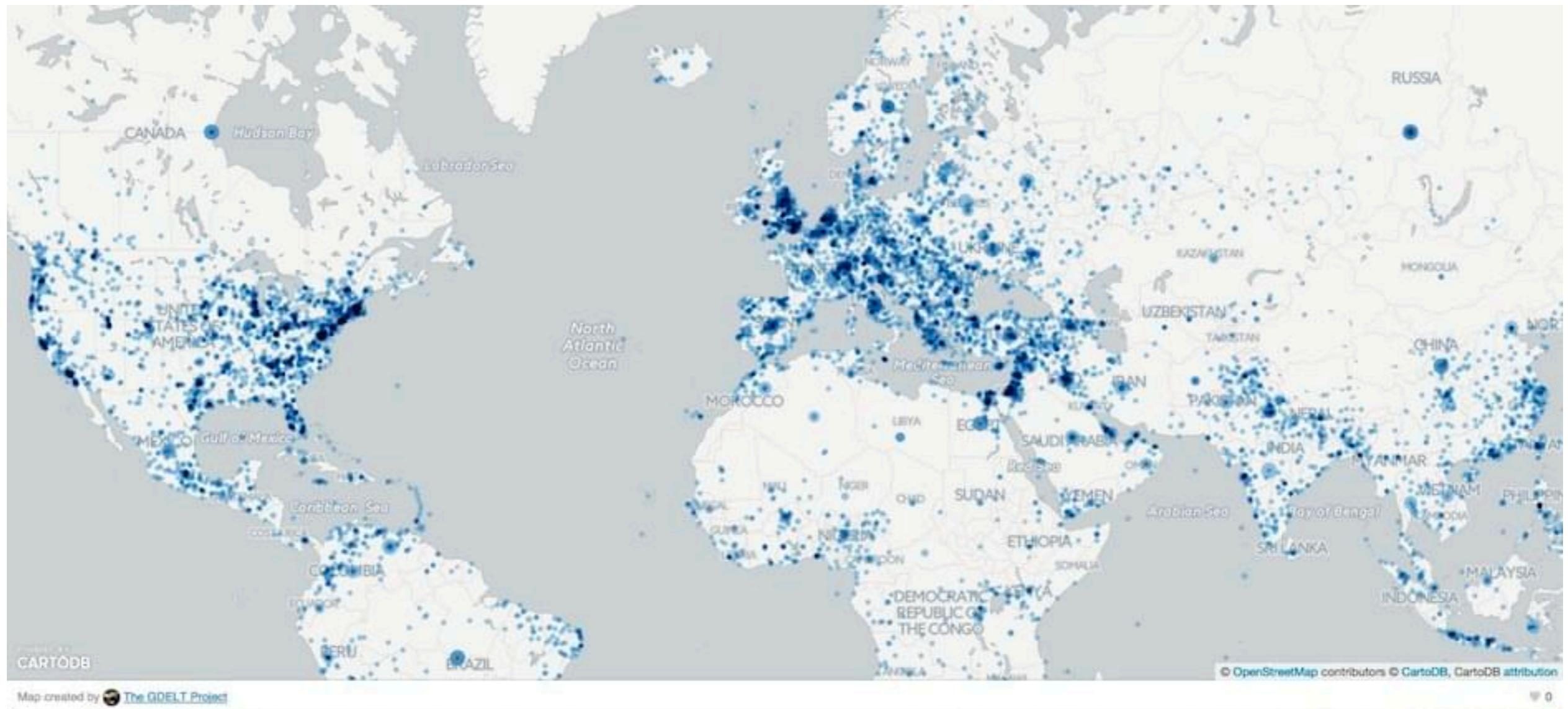
- Architecture EMR :

GDELT

AMÉLIORATIONS



AMÉLIORATIONS POSSIBLES



TELECOM PARISTECH - 2019

GDELT - NOSQL

