

# Projet INF-728

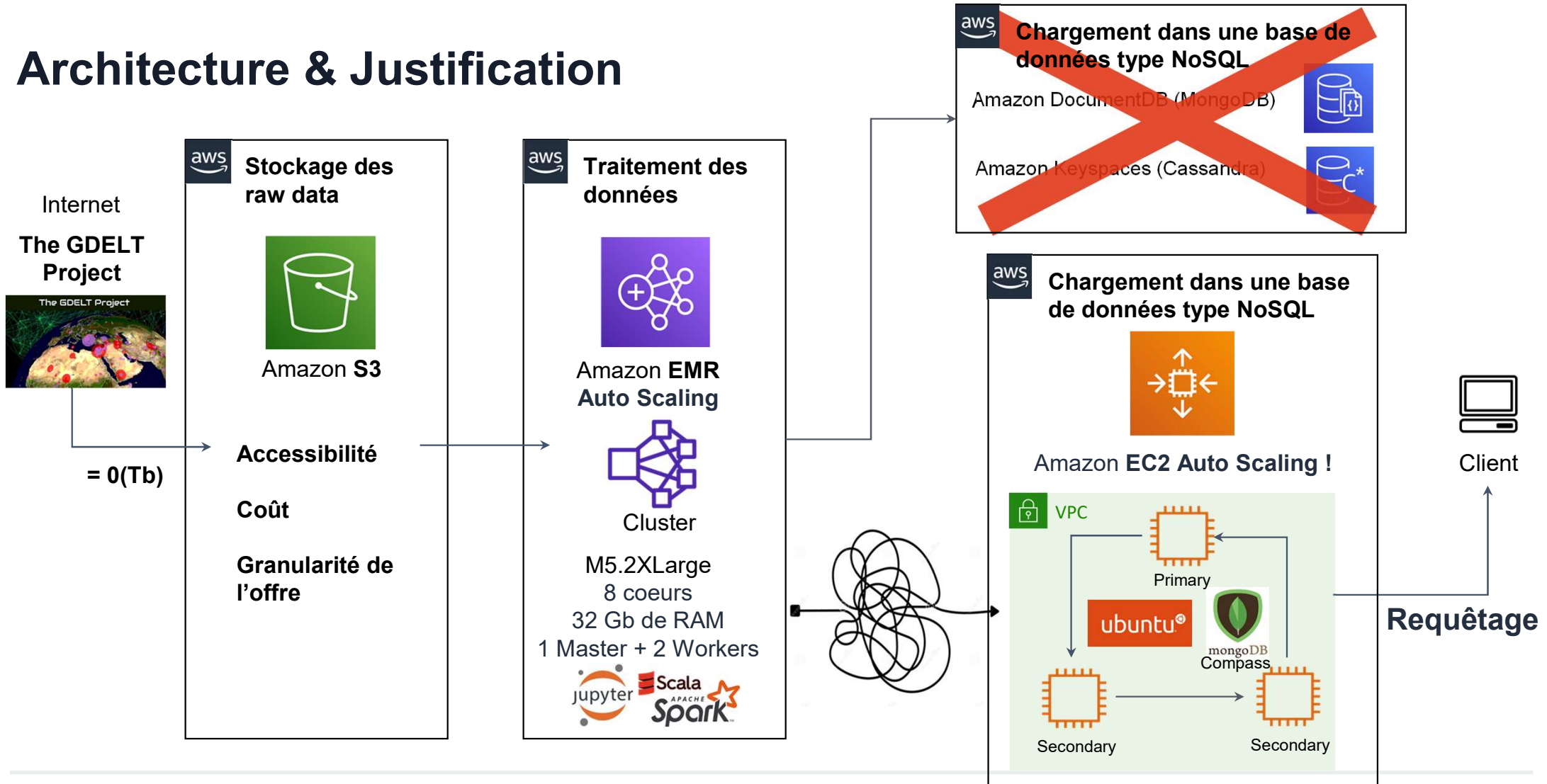
## MS BGD 2020-2021 – 22 jan 2021



Nicolas Calligaro  
Frederic Haykal  
Julien Maksoud  
Axel Michalewicz  
Lingli Zhang



# Architecture & Justification



# Cassandra VS MongoDB

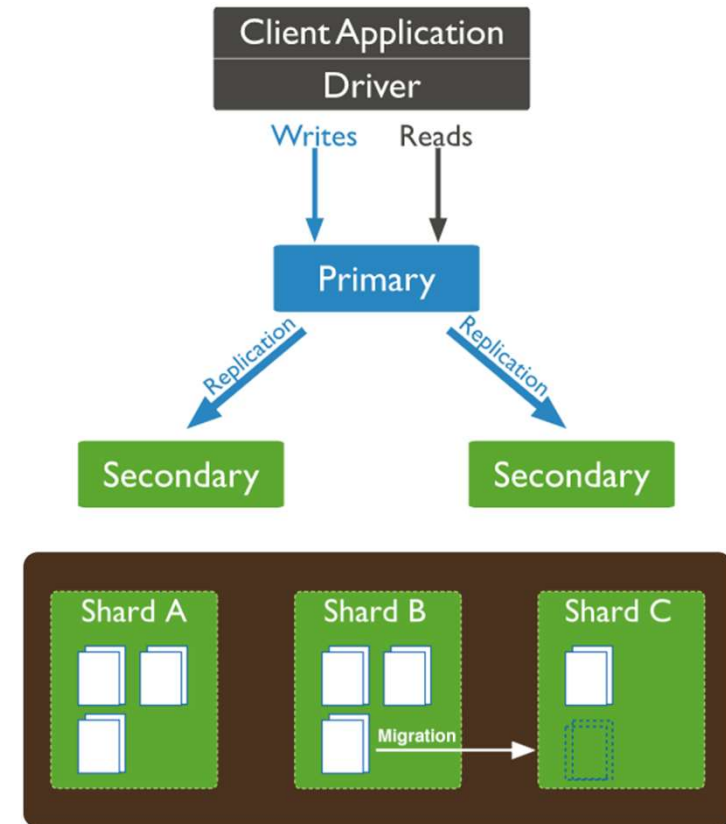
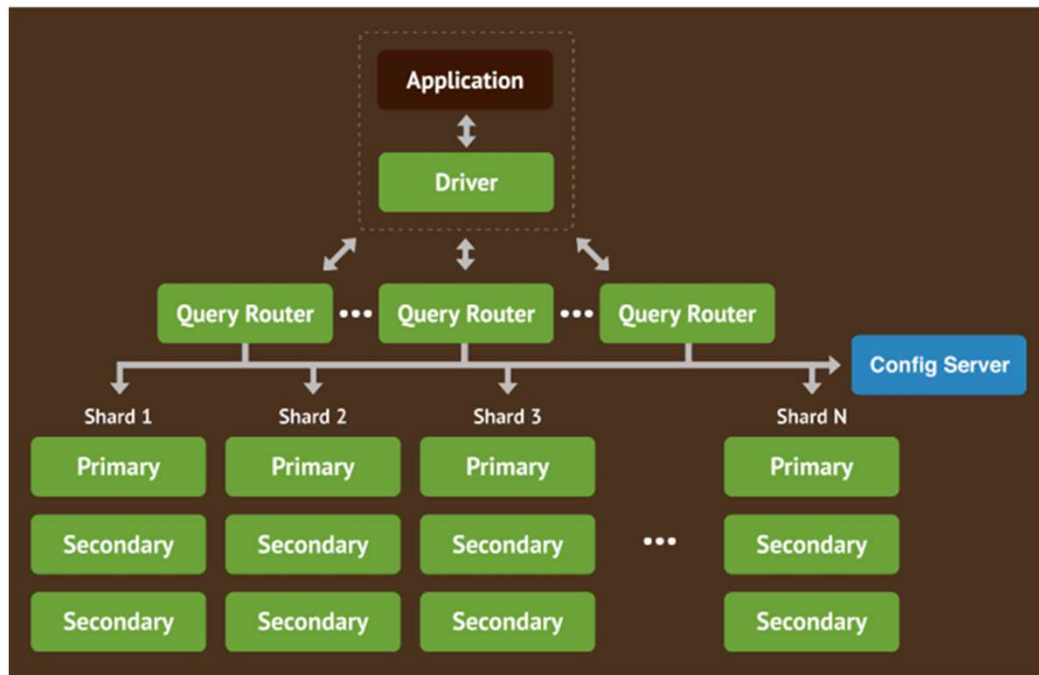
Caractéristique	Important pour ce projet	Cassandra	MongoDB
<i>Modèle De Donnée</i>	Oui	Colonnes	+ Documents
<i>Langage De Requête</i>	Oui	Léger et simple	+ Mongo Query Langage (Puissant)
<i>Agrégation</i>	Oui	Néant	+ Aggregation Framework (Compass)
<i>Nœud Principal</i>	Non	+ Multiple	Unique ou Sharding
<i>Passage À L'échelle</i>	Non	+ Native	Null ou Sharding
<i>Schéma</i>	Non	Fixe	+ Schéma Modifiable à Chaud
<i>Chargement Des Données</i>	Non	+ Efficace	Relative
<i>Consistance Des Donnée</i>	Oui	+ Meilleur	Relative (Réplicat Absorbe Crazy Monkey)

## Besoins du projet GDELT

- Pas de flux continu → Passage à l'échelle pas primordial
- Requêtes complexes de documents nécessaires
- Réduire la quantité de données / Minimiser le nombre de tables



# Architecture MongoDB (cf cours de Geoffrey)



*Load Balancing*

# Import des données dans AWS S3

1 fichier de chaque toutes les 15 minutes: masterfilelist.txt

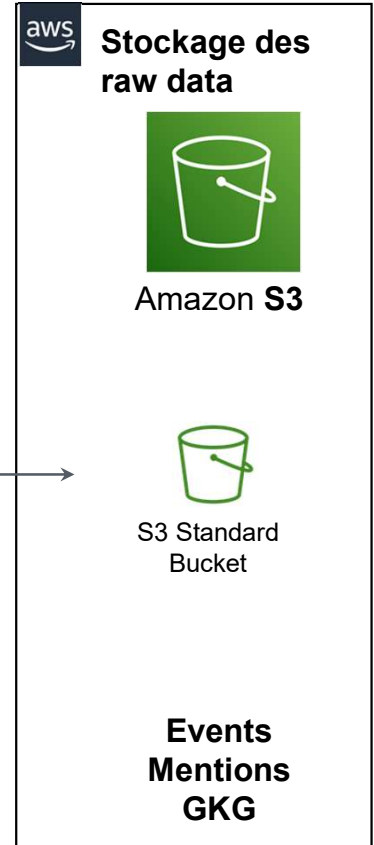
```
masterfilelist.txt
117660 6504649 3e4882e7b5aa79b94b9d800acabdb97 http://data.gdeltproject.org/gdeltv2/20210108120000.gkg.csv.zip
117661 68696 8392b1de9eaf3ce7b7ec2fc70839d767 http://data.gdeltproject.org/gdeltv2/20210108121500.export.CSV.zip
117662 125923 b876c5f98da95bc624e88f98fcb1af62 http://data.gdeltproject.org/gdeltv2/20210108121500.mentions.CSV.zip
117663 4748976 9cbfa559aab6f13f27511d02dalaafd00 http://data.gdeltproject.org/gdeltv2/20210108121500.gkg.csv.zip
117664 92211 643075007162c647026540d136317b4b http://data.gdeltproject.org/gdeltv2/20210108123000.export.CSV.zip
117665 172181 a846ebfb98523423a8f6d658778e721 http://data.gdeltproject.org/gdeltv2/20210108123000.mentions.CSV.zip
117666 6265775 b3ad64294216ab836e8a0a3dd6647bb3 http://data.gdeltproject.org/gdeltv2/20210108123000.gkg.csv.zip
117667 86010 e04ec2ba8096a5dd91fe83e60fed022 http://data.gdeltproject.org/gdeltv2/20210108124500.export.CSV.zip
117668 158964 fd818796b9198aa281ac01c72bdf27 http://data.gdeltproject.org/gdeltv2/20210108124500.mentions.CSV.zip
117669 6310064 5ac32dd4e374b948579407697cee6 http://data.gdeltproject.org/gdeltv2/20210108124500.gkg.csv.zip
117670 87435 2ff9a44b8efa604a7400125f6daa683 http://data.gdeltproject.org/gdeltv2/20210108124500.export.CSV.zip
117671 164973 84045cac0c96ae26d75e8c34add649 http://data.gdeltproject.org/gdeltv2/20210108124500.mentions.CSV.zip
117672 6133097 4cc4d52d4070081536665168cf24d http://data.gdeltproject.org/gdeltv2/20210108124500.gkg.csv.zip
117673 73182 9fdb629b84f04bdde7908a488147d8 http://data.gdeltproject.org/gdeltv2/20210108124500.export.CSV.zip
117674 161068 3aaf3d3d90e0691f94d4f22204a16bb http://data.gdeltproject.org/gdeltv2/20210108124500.mentions.CSV.zip
117675 6873515 c6c0313c4af0cd12eb8e11fcd5773 http://data.gdeltproject.org/gdeltv2/20210108124500.gkg.csv.zip
117676 85725 1c8bcccc8c92c81645ed55b986d360e http://data.gdeltproject.org/gdeltv2/20210108124500.export.CSV.zip
117677 162038 034e04dc8d9a0f4185fe6b8ccf9c3a http://data.gdeltproject.org/gdeltv2/20210108124500.mentions.CSV.zip
117678 6705445 e26bd8aa05d79c5377c47cf3ee468 http://data.gdeltproject.org/gdeltv2/20210108124500.gkg.csv.zip
117679 80973 179f442631181bd2160e4d0be6413ac http://data.gdeltproject.org/gdeltv2/20210108124500.export.CSV.zip
117680 172948 b485ef1c3d132174173b65149edbed http://data.gdeltproject.org/gdeltv2/20210108124500.mentions.CSV.zip
```

Lancement du téléchargement des fichiers de GDELT vers S3

```
Entrée [13]: 1 def StokeFileS3(URL: String) = {
2             2 val fileName = URL.split("/")
3             3 val dir = "/mnt/tmp/"
4             4 val localFileName = dir + fileName
5             5 try {
6               fileDownloader(URL, localFileName)
7               @transient val s3Client: AmazonS3 = AmazonS3ClientBuilder.defaultClient()
8               val localFile = new File(localFileName)
9               s3Client.putObject("testfuret/Master_Translate_file ", fileName, localFile)
10              localFile.delete()
11            } catch {
12              case e: java.io.FileNotFoundException =>
13              case e: com.amazonaws.sdk.ClientException =>
14              case e: java.io.IOException =>
15              case e: java.lang.RuntimeException =>
16              case e: java.lang.Exception =>
17            }
18          }
```

```
size hash url
6178636|72d0e5ca18362beca19a0510854f4721|http://data.gdeltproject.org/gdeltv2/20200201000000.translation.gkg.csv.zip
6133442|765b9da4dd3270732defd66de74be37|http://data.gdeltproject.org/gdeltv2/20200201001500.translation.gkg.csv.zip
6454334|209299355a96b765a9ecf3ec77d2f6de|http://data.gdeltproject.org/gdeltv2/20200201003000.translation.gkg.csv.zip
5654513|1dd9159fca3a3ebf01be064bdc4bfa|http://data.gdeltproject.org/gdeltv2/20200201004500.translation.gkg.csv.zip
5517145|5570a7f6cf093ec26464d07d8d49846|http://data.gdeltproject.org/gdeltv2/20200201006000.translation.gkg.csv.zip
5602121|2c41db3dc8db82d8f0f0d034baca89bd|http://data.gdeltproject.org/gdeltv2/20200201011500.translation.gkg.csv.zip
6049830|22430c0b0f433ec91e4fcb777d1f431e|http://data.gdeltproject.org/gdeltv2/20200201013000.translation.gkg.csv.zip
6369154|39a46116c40134a5db96beeb30f434c|http://data.gdeltproject.org/gdeltv2/20200201014500.translation.gkg.csv.zip
5541150|cb02ef523a91f05a31f1f519f8b0516|http://data.gdeltproject.org/gdeltv2/20200201020000.translation.gkg.csv.zip
4966893|8d8112546a147bc3a9d1322176ab3f0e|http://data.gdeltproject.org/gdeltv2/20200201021500.translation.gkg.csv.zip
5464143|35e3848d0879f09c686a6aaf086f|http://data.gdeltproject.org/gdeltv2/20200201023000.translation.gkg.csv.zip
5997714|7567655905c804f9819f9e86e2796|http://data.gdeltproject.org/gdeltv2/20200201024500.translation.gkg.csv.zip
5524668|e6d69c9b098990269c2a975f1cad|http://data.gdeltproject.org/gdeltv2/20200201030000.translation.gkg.csv.zip
5202802|44523f0855f67a40d5419b3f9963ee|http://data.gdeltproject.org/gdeltv2/20200201031500.translation.gkg.csv.zip
4663946|3e06f6d1faa36a72eeefaf3ad03a9a7|http://data.gdeltproject.org/gdeltv2/20200201033000.translation.gkg.csv.zip
6223891|0680af7745c352c8caeeb8f3270fdd|http://data.gdeltproject.org/gdeltv2/20200201034500.translation.gkg.csv.zip
```

Stockage des raw data (S3)  
Traitement (EMR)  
Mise en DB NoSQL (EC2)  
Requêtage (EC2)



Quelques Tb en CSV



# Traitement des données dans EMR

Stockage des raw data (S3)  
Traitement (EMR)  
Mise en DB NoSQL (EC2)  
Requêtage (EC2)

