

PROJET GDELT

How was 2021 ?

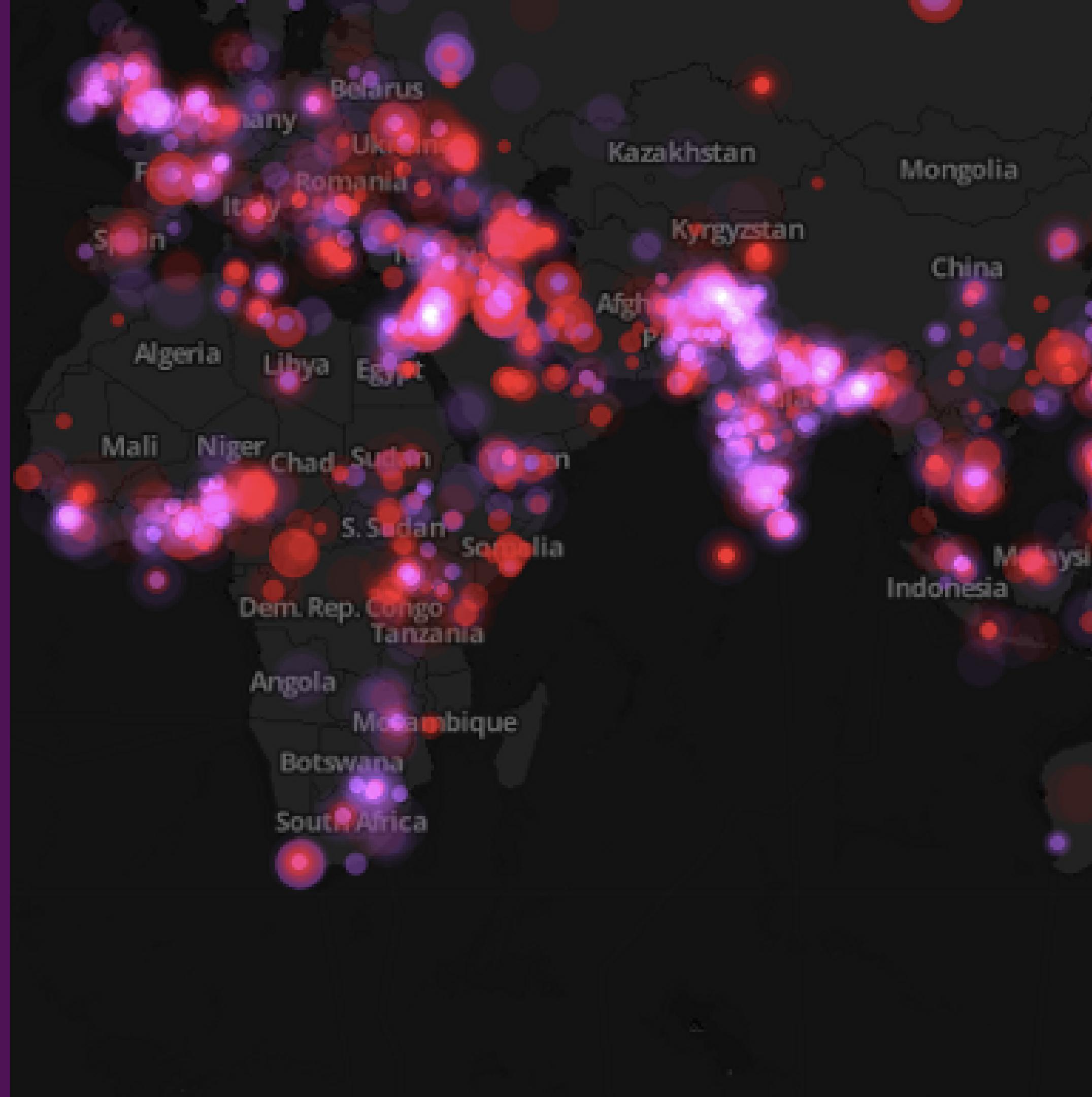
Sommaire

1. Contexte
2. Description des données
3. Architecture
4. Requêtage
5. Demo
6. Améliorations

Contexte

Objectif:

- Analyser l'impact media de l'année 2021 via le jeu des données GDELT.
- Proposer un système de stockage distribué, résilient et performant pour les données GDELT.



GDELT Dataset

Events

Ce dataset regroupe les informations des événements en capturant les 2 participants à cet événement.

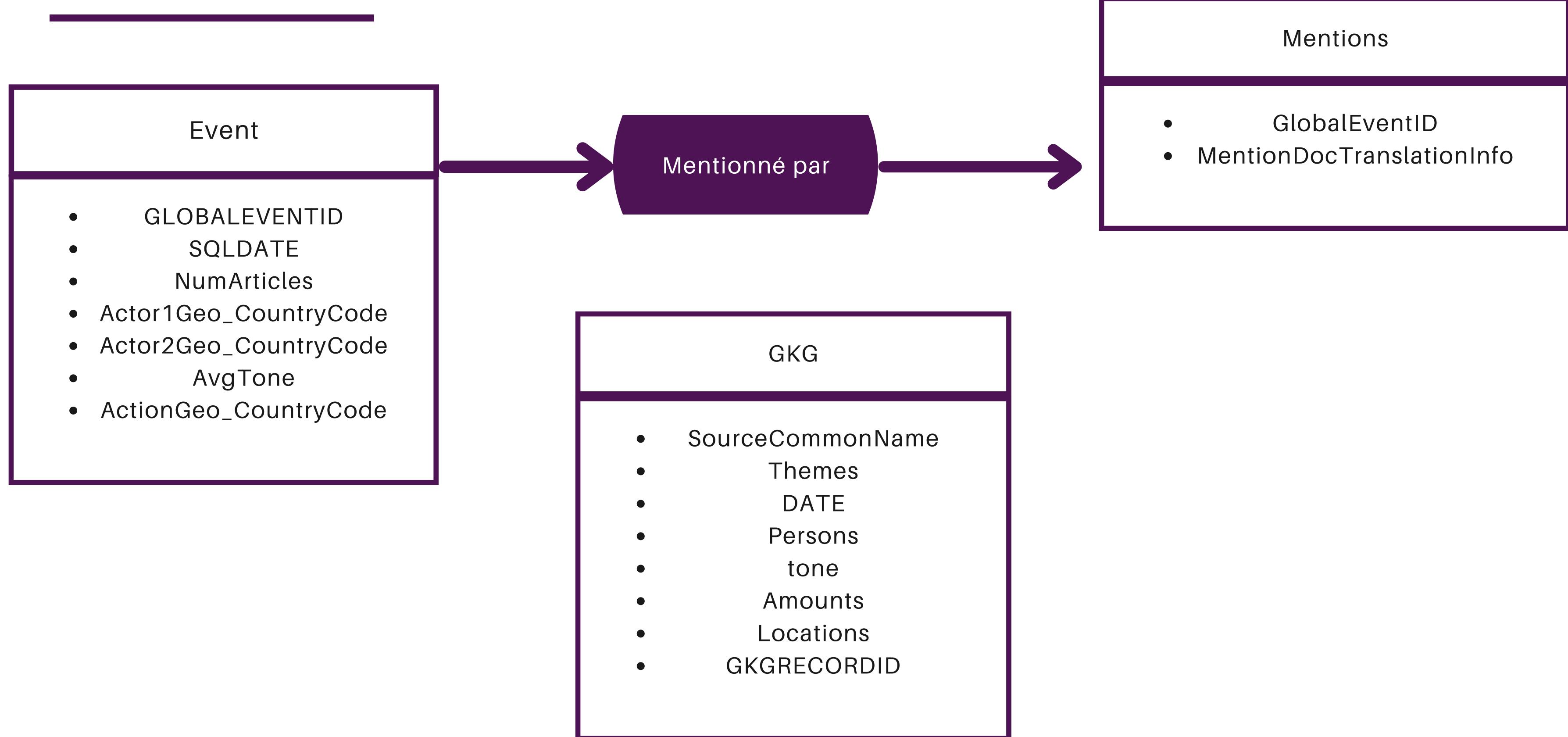
Mentions

Cette table regroupe toutes les mentions pour un événement .

GKG

Global Knowledge Graph permet de connecter les participants, organisations, lieux et thèmes.

Modèle des données



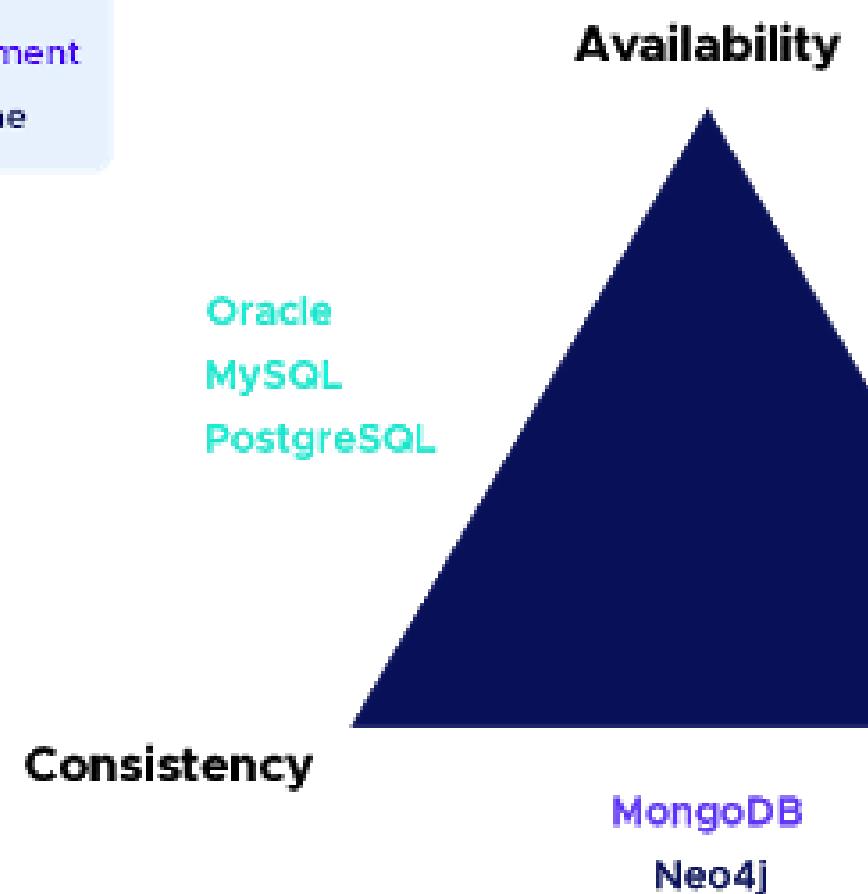
Choix de la technologie:

MongoDB:

- Fonctionnalité.
- Cohérence
- Tolérance au partitionnement
- Sharding - Multi-indexing.
- Rapidité de requêtage



■ Relationnel
■ Orienté Document
■ Orienté Graphe



Gdelt



Extraction des données

Luigi

Stockage des données + requêtes



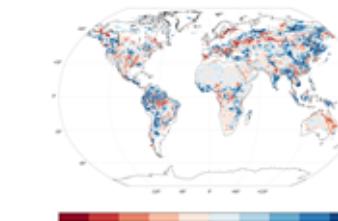
6 machines

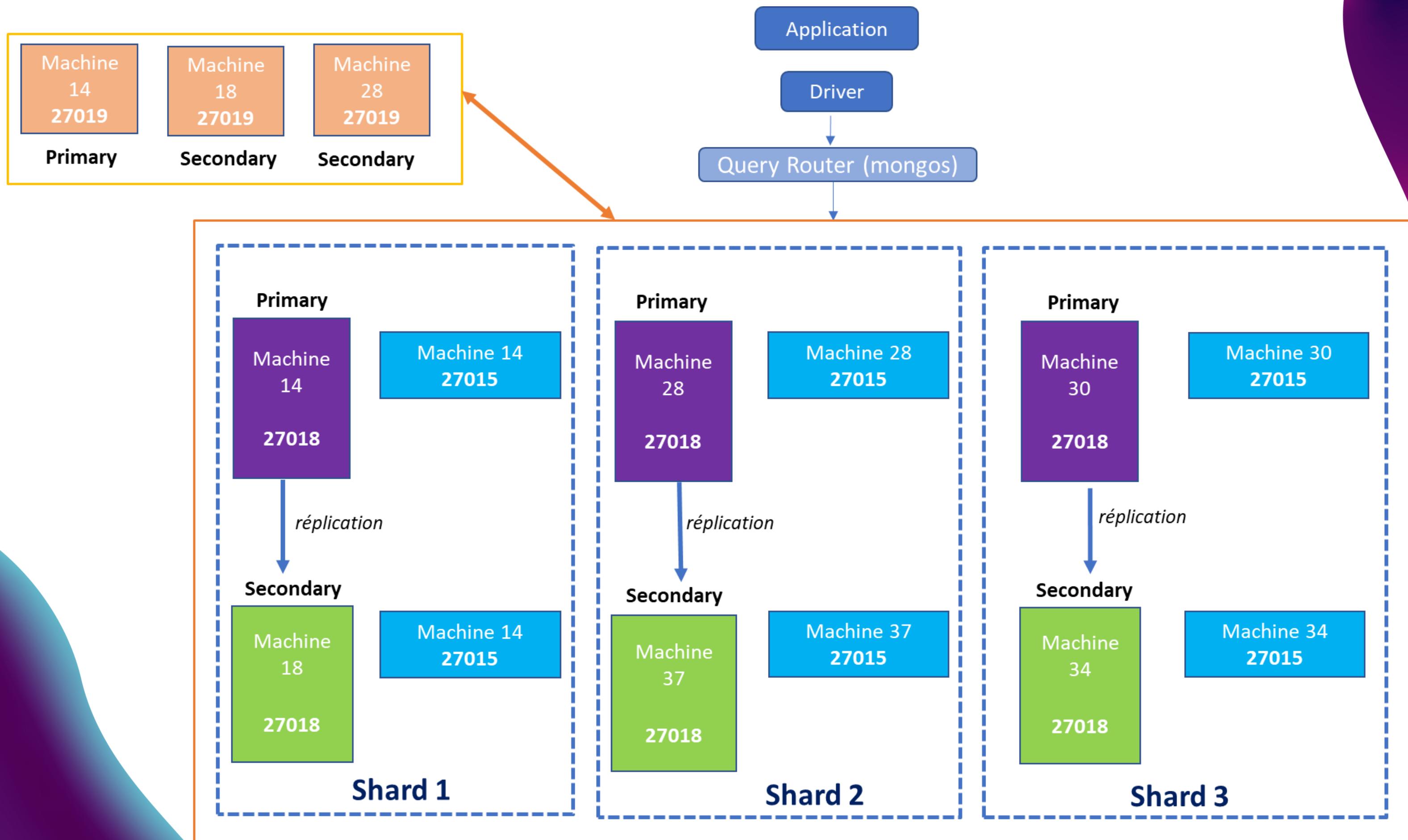


mongoDB



Requêtes + résultats + analyses





ETL

Luigi



```
import luigi

class MyTask(luigi.Task):
    param = luigi.Parameter(default=42)

    def requires(self):
        return SomeOtherTask(self.param)

    def run(self):
        f = self.output().open('w')
        print >>f, "hello, world"
        f.close()

    def output(self):
        return luigi.LocalTarget('/tmp/foo/bar-%s.txt' % self.param)

if __name__ == '__main__':
    luigi.run()
```

The business logic of the task

Where it writes output

What other tasks it depends on

Parameters for this task

Luigi Task Status

Task List Dependency Graph Workers Resources Running

TASK FAMILIES Clear selection

	Name	Details	Priority	Time	Actions
14	PartsReport.CombineReports				
44	PartsReport.DatedBOReport				
16	PartsReport.DownloadReport				
16	PartsReport.ExtractReport				
2	PartsReport.GenerateFiles				
6	PartsReport.LoadTable	table=FUP, date=2019-05-28	0	5/28/2019, 8:00:17 AM	
2	PartsReport.LoadPartPlant				
2	PartsReport.LoadPartPlantUse				
2	PartsReport.LoadParts				
2	PartsReport.DebateScrum				
	PFEP				
	BillOfMaterials				
	Others				
1	LoadTable				

PENDING TASKS: 3 RUNNING TASKS: 1 BATCH RUNNING TASKS: 0 DONE TASKS: 134

FAILED TASKS: 3 UPSTREAM FAILURE: 1 DISABLED TASKS: 0 UPSTREAM DISABLED: 0

Displaying tasks of family LoadTable.

Show 10 entries Filter table: Filter on Server

Showing 1 to 1 of 1 entries (filtered from 142 total entries) Previous 1 Next

The screenshot displays the Luigi Task Status interface. On the left, a sidebar lists task families: data_management, PartsReport, PFEP, BillOfMaterials, Others, and a specific LoadTable entry. The main area shows summary statistics for pending, running, batch running, and done tasks, along with failed tasks and upstream failures. Below this is a detailed table of running tasks, with one row for LoadTable. The table includes columns for Name, Details, Priority, Time, and Actions. A message at the bottom indicates it's displaying tasks of the LoadTable family.

STACKTRACE (most recent call first)

[Full](#) [Raw](#)

```
__init__.py in handle at line 894
889.     """
890.     rv = self.filter(record)
891.     if rv:
892.         self.acquire()
893.         try:
894.             self.emit(record)
895.         finally:
896.             self.release()
897.     return rv
898.
899. def setFormatter(self, fmt):
```

record	<LogRecord: xxx>
rv	True
self	<EventHandler (WARNING)>

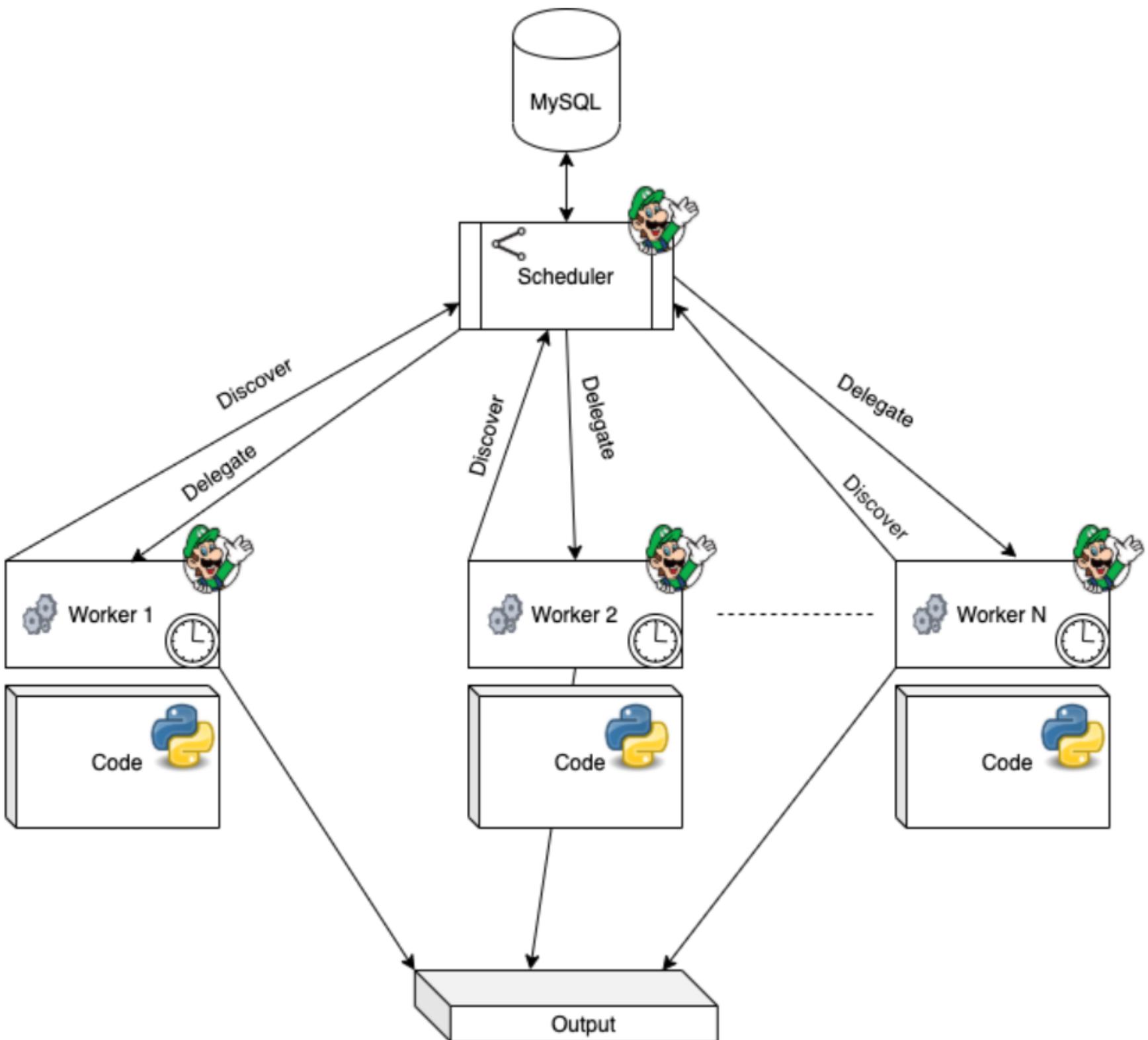
```
__init__.py in handle at line 1524
```

```
__init__.py in _log at line 1514
```

```
__init__.py in warning at line 1390
```

```
xxx/apiclient.py in _get at line 52
```

```
47.         'x-api-key': self.api_key,
48.     }
49.
50.     def _get(self, url, **kwargs):
51.         # Log warning
52.         log.warning('Test warning', extra={'test': 'value'})
```

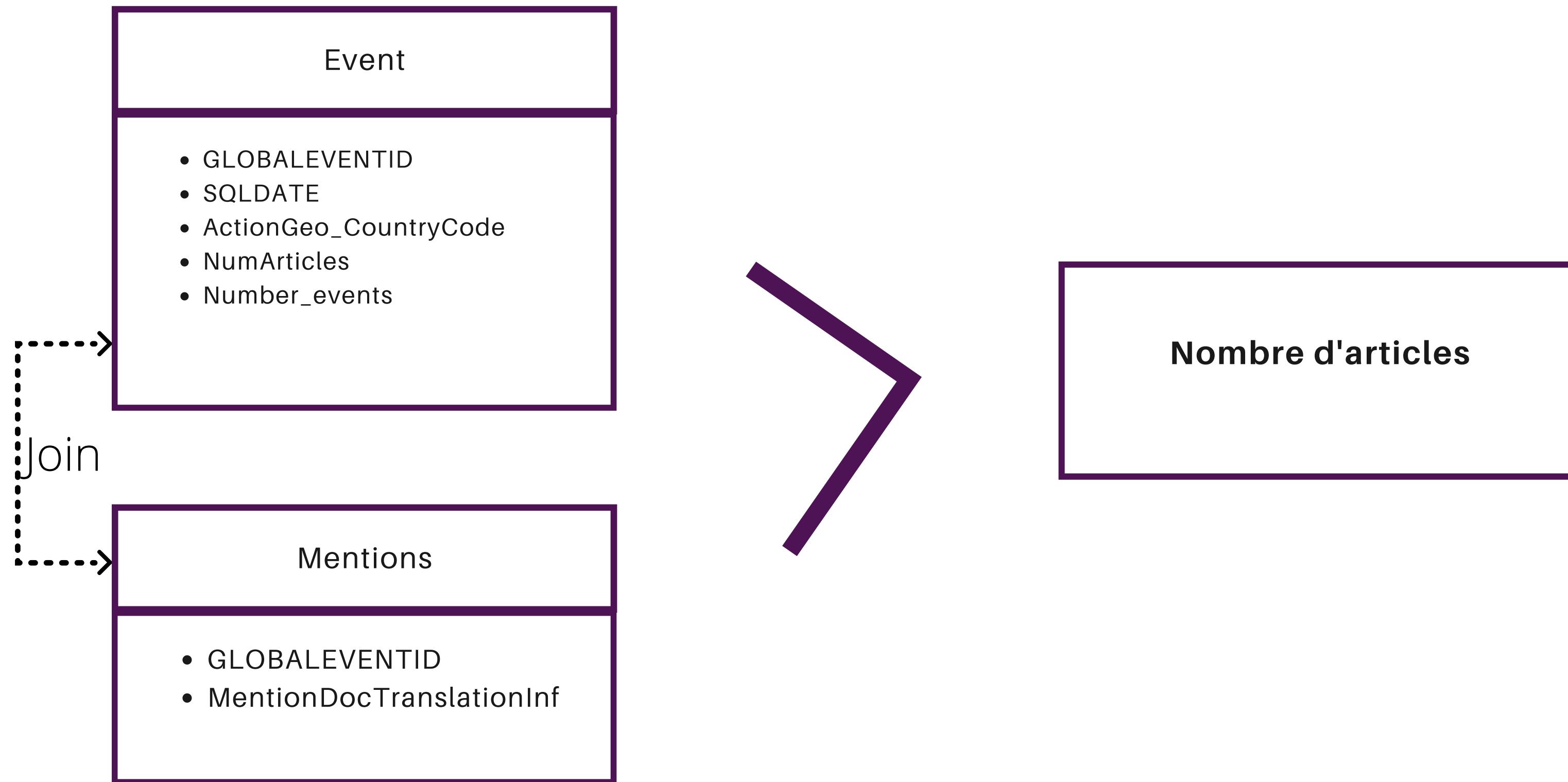


Luigi workers communication with the scheduler



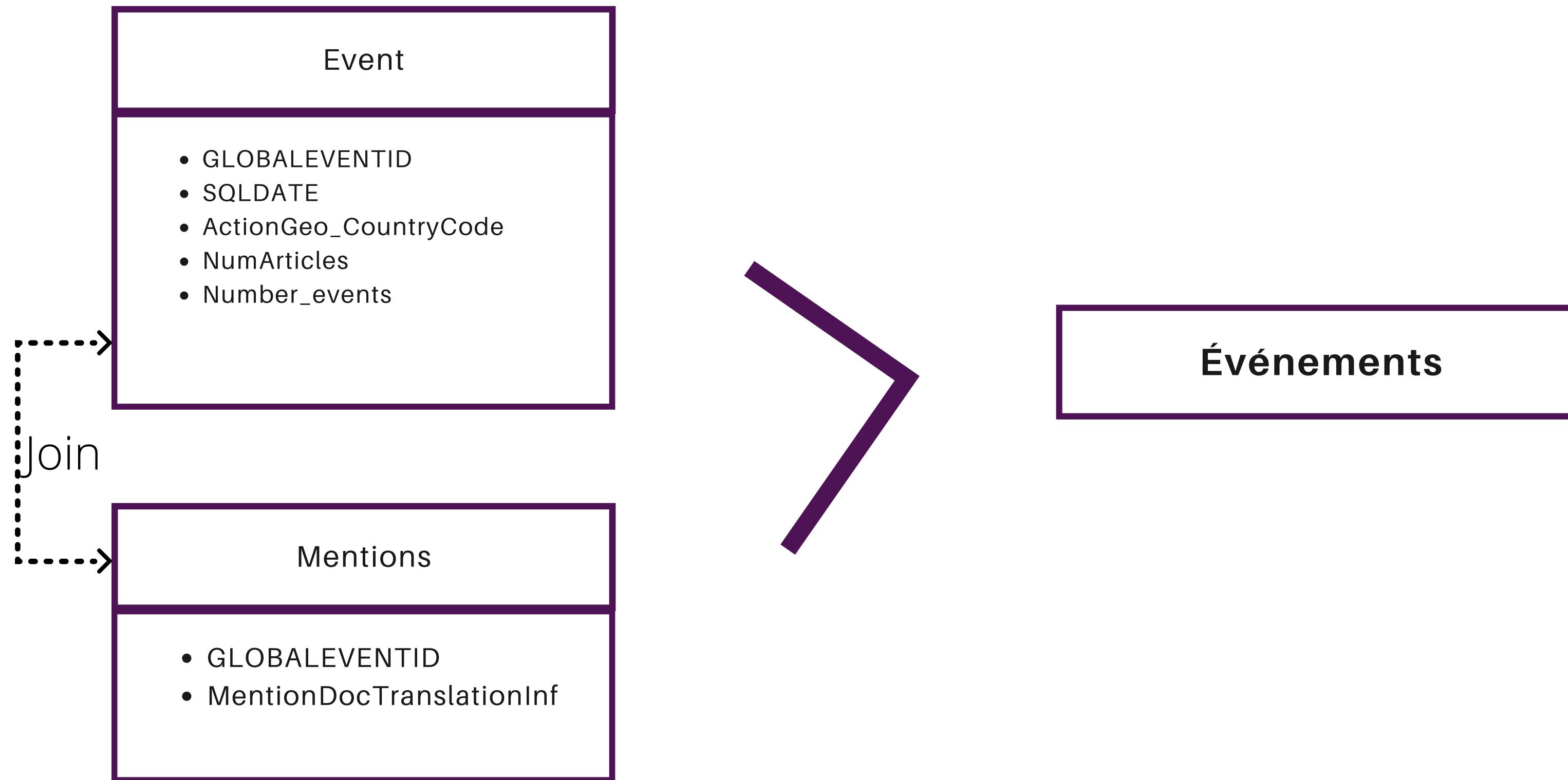
Requête 1 :

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



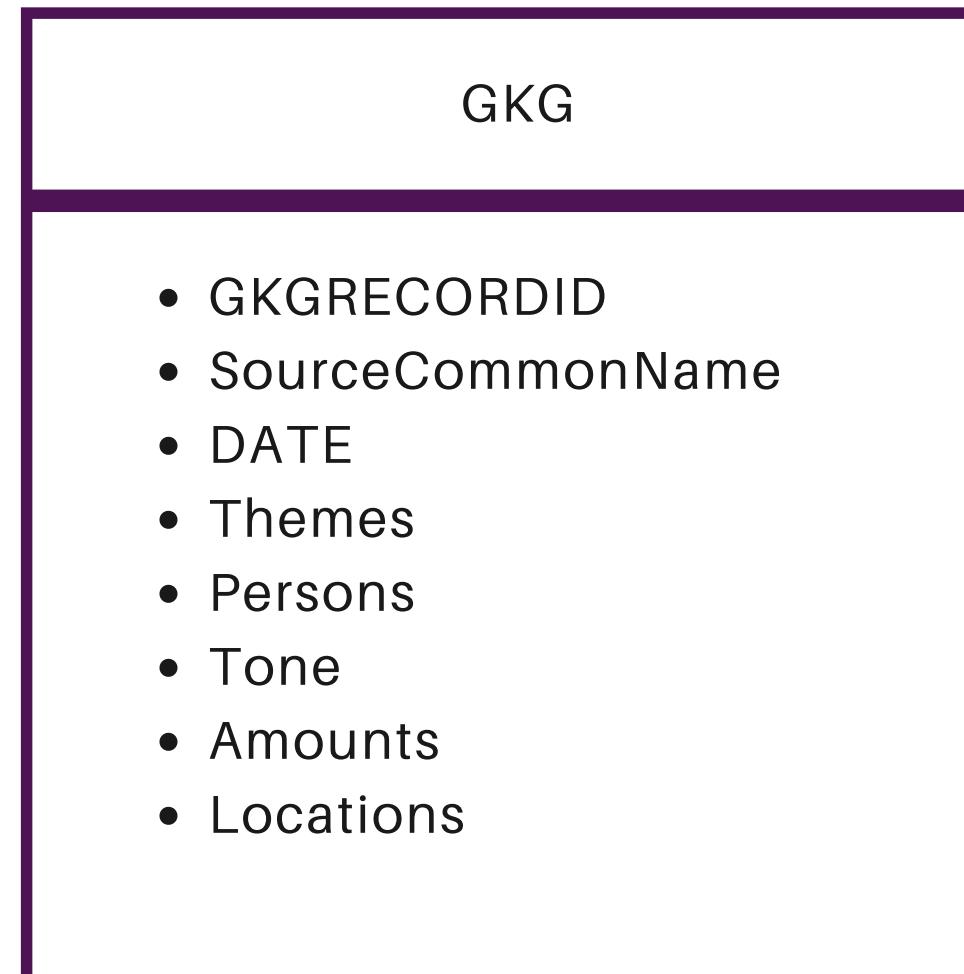
Requête 2 :

Pour un pays donné en paramètre, affichez les évènements qui y ont eu place triées par le nombre de mentions (tri décroissant); permettez une agrégation par jour/mois/année



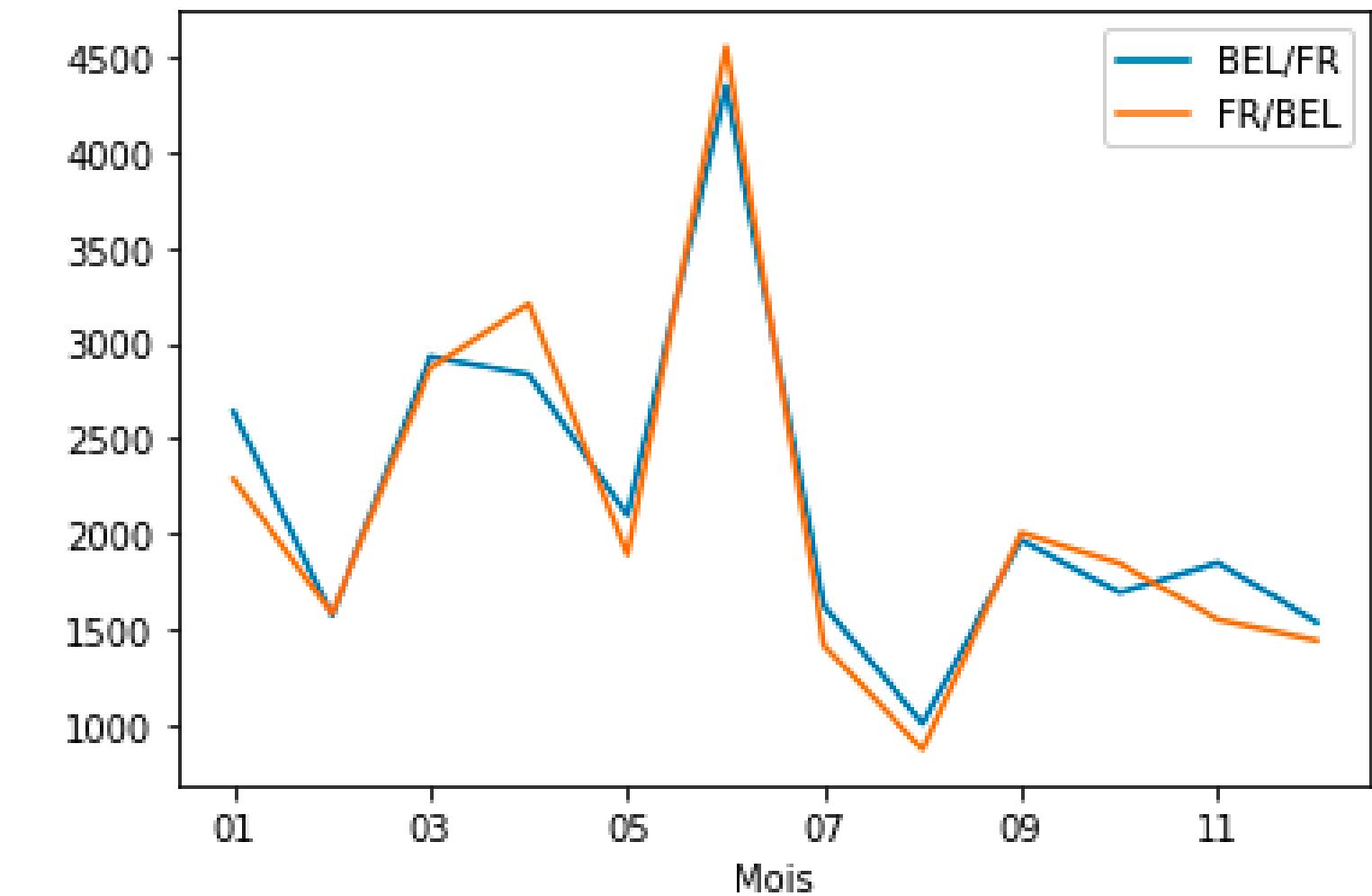
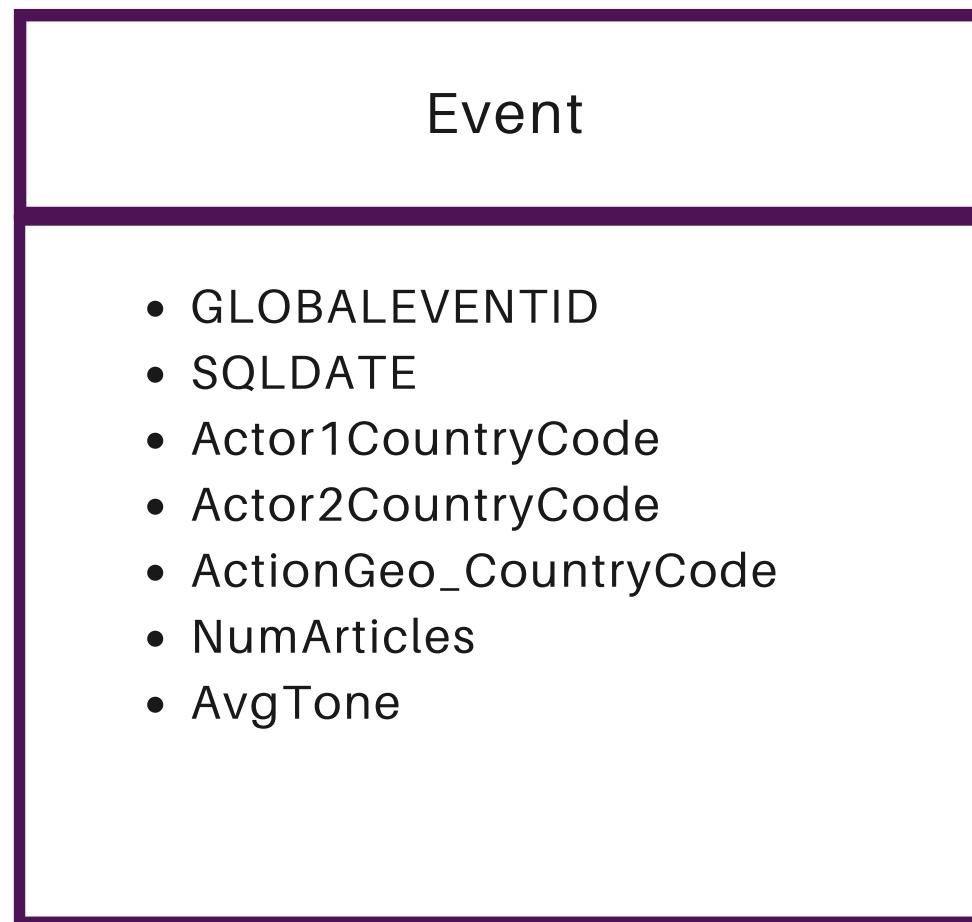
Requête 3 :

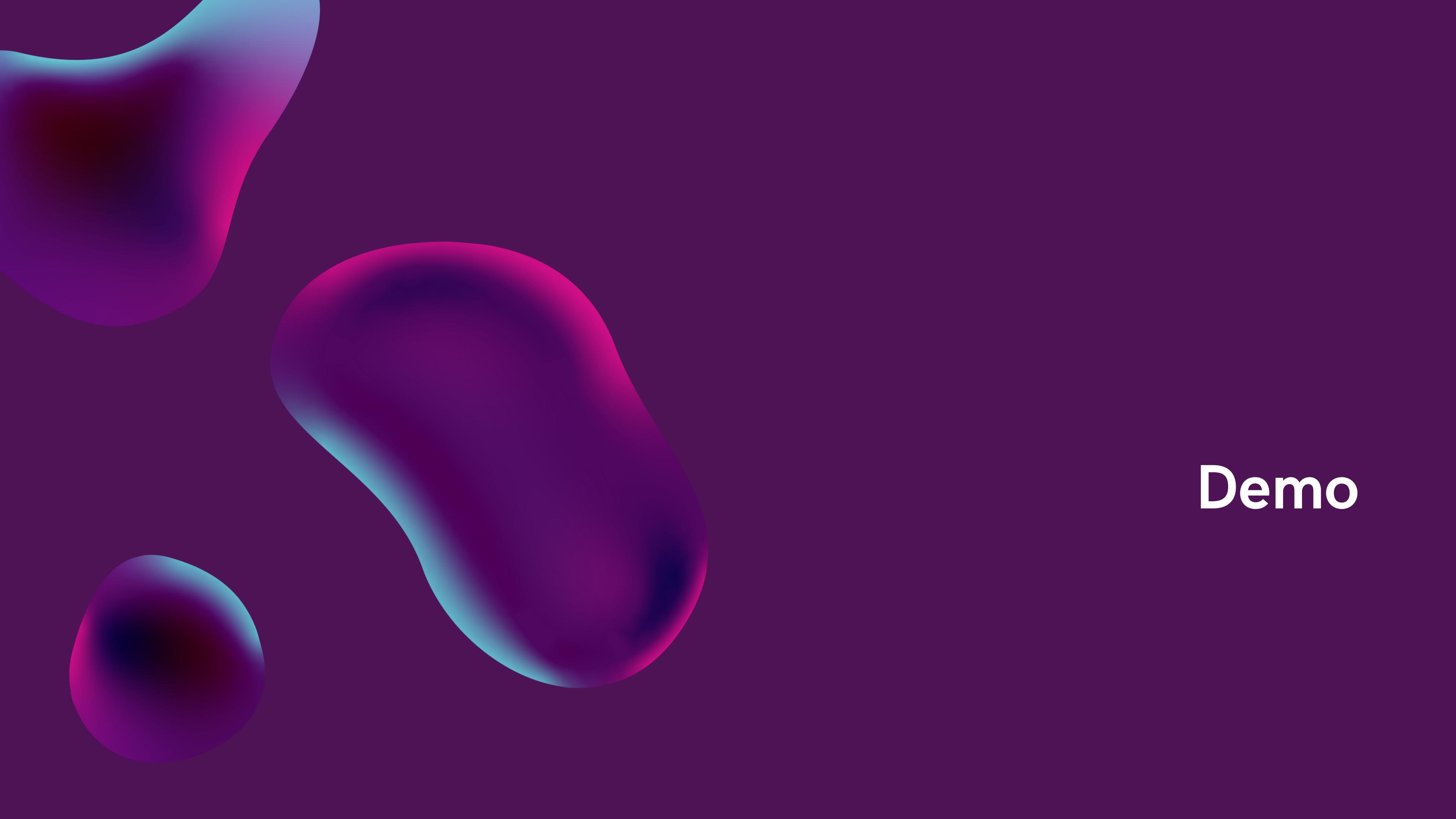
Pour une source de données passée en paramètre (gkg.SourceCommonName) affichez les thèmes, personnes, lieux dont les articles de cette sources parlent ainsi que le nombre d'articles et le ton moyen des articles (pour chaque thème/personne/lieu); permettez une agrégation par jour/mois/année.



Requête 4 :

L'évolution des relations entre deux pays (spécifiés en paramètre) au cours de l'année. Vous pouvez vous baser sur la langue de l'article, le ton moyen des articles, les thèmes plus souvent cités, les personnalités ou tout élément qui vous semble pertinent.



The background features three large, semi-transparent circles with a gradient from dark purple to bright blue. One circle is positioned at the top left, another at the bottom left, and a third, larger one is centered in the middle-left area.

Demo