



# **ANALYSE TEMPS RÉEL DES DISCUSSIONS REDDIT CAN 2025**

## **Architecture Big Data & Machine Learning**

### **Présentée par:**

BOUTANFIT Salma  
ELIDRISSI Asma  
JENNANE Salma  
Ouahib Salma

### **Jury:**

Yasser EL MADANI  
EL ALAMI

# Plan

- 1 Contexte & Problématique**
- 2 Objectifs du projet**
- 3 Architecture du système**
- 4 Réalisation**
- 5 Difficultés rencontrés et solutions**
- 6 Conclusion et perspectives**



# 1 Contexte & problématique



# 1.Contexte



- Événement sportif majeur africain
- Millions d'interactions sur réseaux sociaux



- Plateforme de discussions thématiques
- Communautés (subreddits) engagées
- Source riche d'opinions et analyses

## DÉFIS BIG DATA



**Volume massif** • Centaines de posts par heure



**Flux continu** • Données en temps réel



**Données non structurées** • Texte libre, langues multiples



**Traitement temps réel** • Analyses instantanées

## QUESTION CENTRALE

- **Comment concevoir une architecture *Big Data* capable de collecter, traiter et analyser en temps quasi-réel les discussions Reddit sur la CAN 2025 ?**

# 2 Objectifs du projet



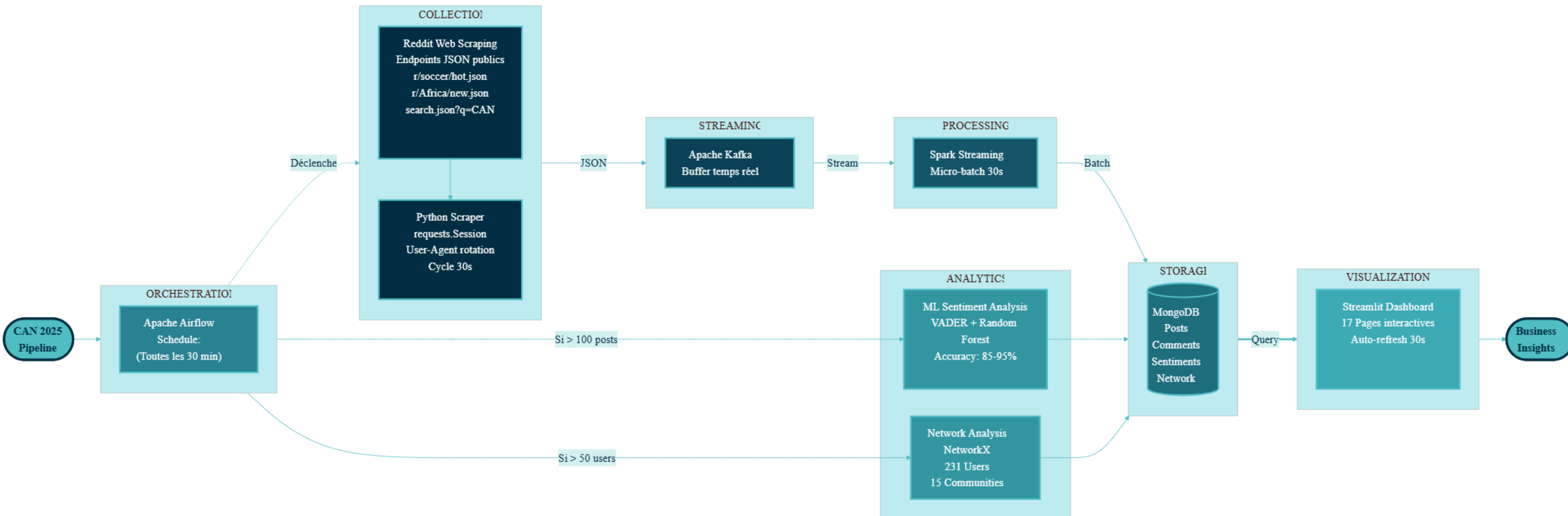
# Objectifs du projet

- ✓ PIPELINE TEMPS RÉEL
- ✓ ANALYSE DE SENTIMENTS (ML)
- ✓ ANALYSE DE RÉSEAU SOCIAL
- ✓ DASHBOARD INTERACTIF

# 4 Architecture du système







# POURQUOI KAPPA ?

## ✓ **SIMPLICITÉ**

Une seule couche de traitement

## ✓ **LATENCE FAIBLE**

Données affichées en < 2 minutes

## ✓ **REJOUABILITÉ**

Kafka permet retraitement historique

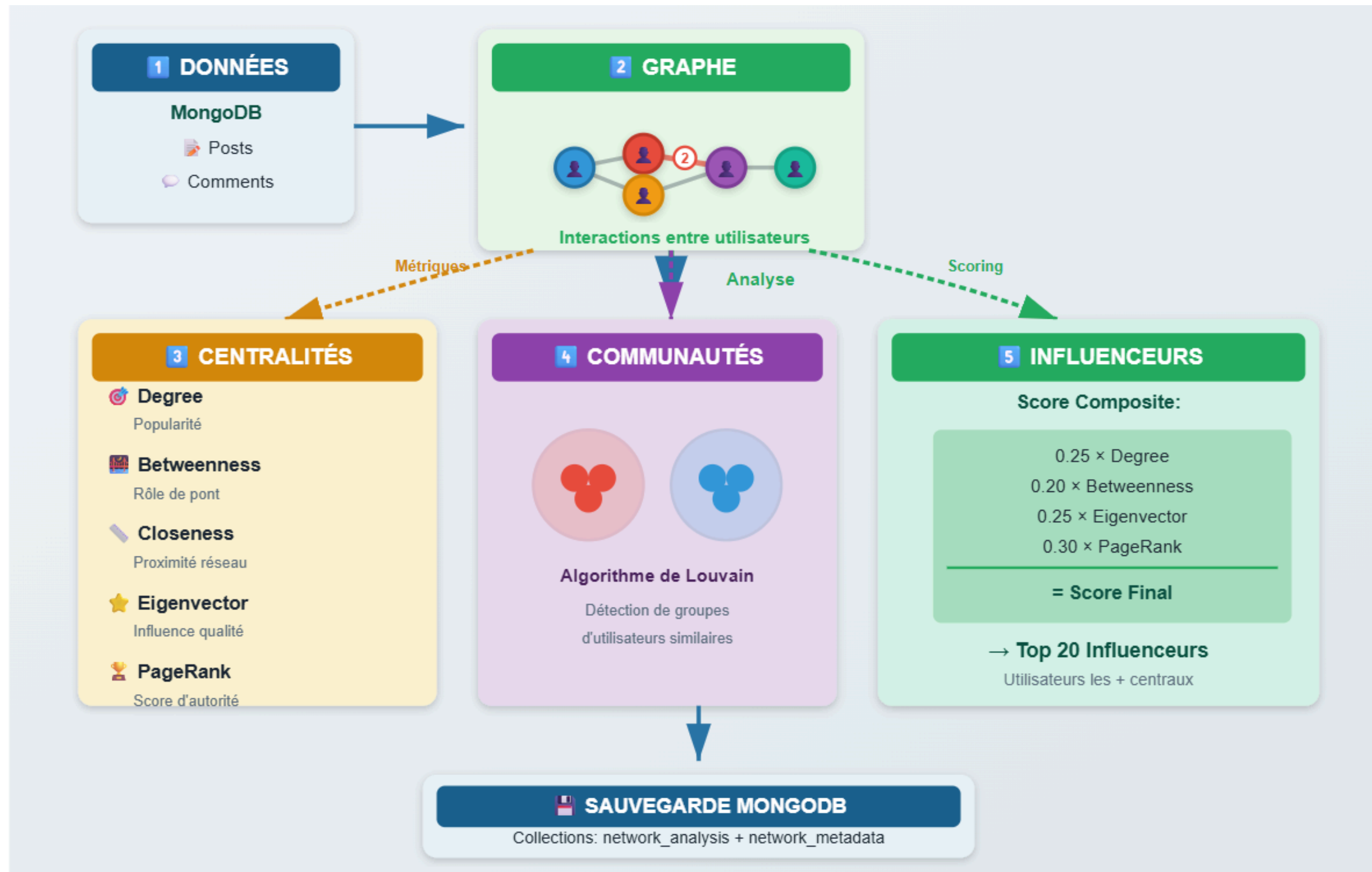
## ✓ **SCALABILITÉ**

Ajout de consumers facilité

Composant	Latence	Fonction
Scraping	30 s	Collecte Reddit
Kafka	< 1 s	Message broker
Spark	30 s	Traitement
MongoDB	Instantané	Écriture
ML Analysis	Variable	Périodique
Network Analysis	Variable	Périodique
Dashboard	30 s	Refresh cache

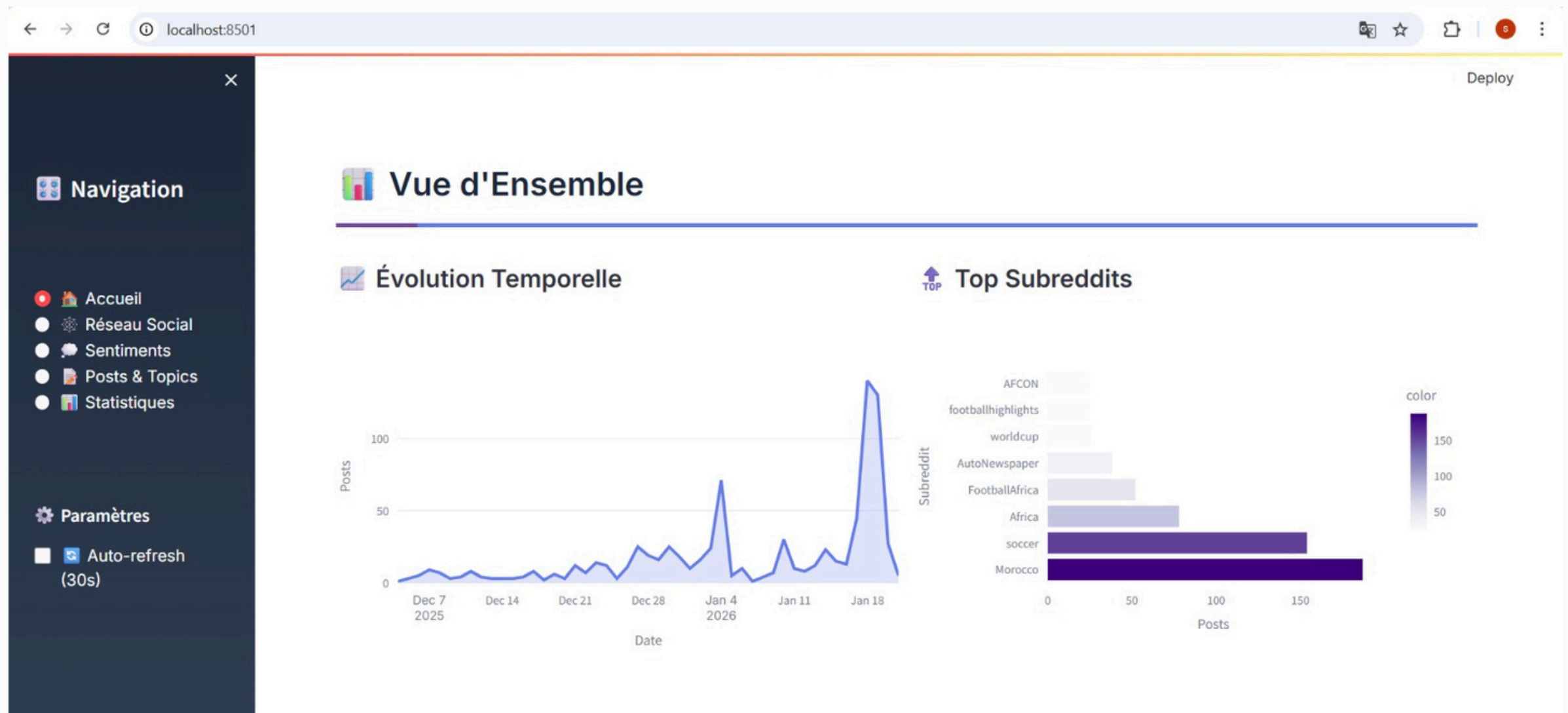
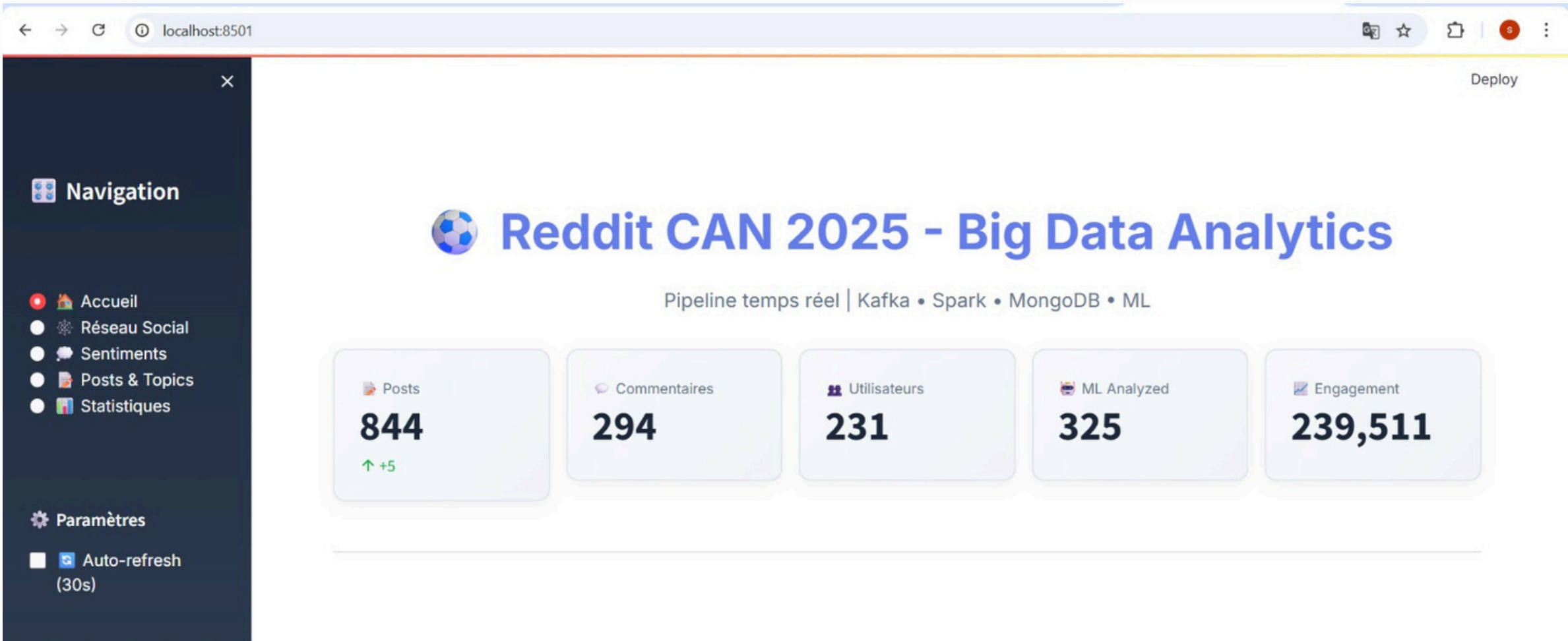
Métrique	Valeur
<b>THROUGHPUT</b>	
Posts/batch Spark	30 posts
Capacité théorique	1,000+ posts/heure
Posts collectés (réels)	844 posts
<b>DISPONIBILITÉ</b>	
Pipeline	24/7
Healthchecks	Docker
Retry automatique	3 tentatives
Monitoring	Airflow centralisé
<b>STOCKAGE</b>	
Taille moyenne	~1.2 MB / 100 posts
Collections MongoDB	5 collections
Backup	Quotidien possible

# Analyse de Réseau Social

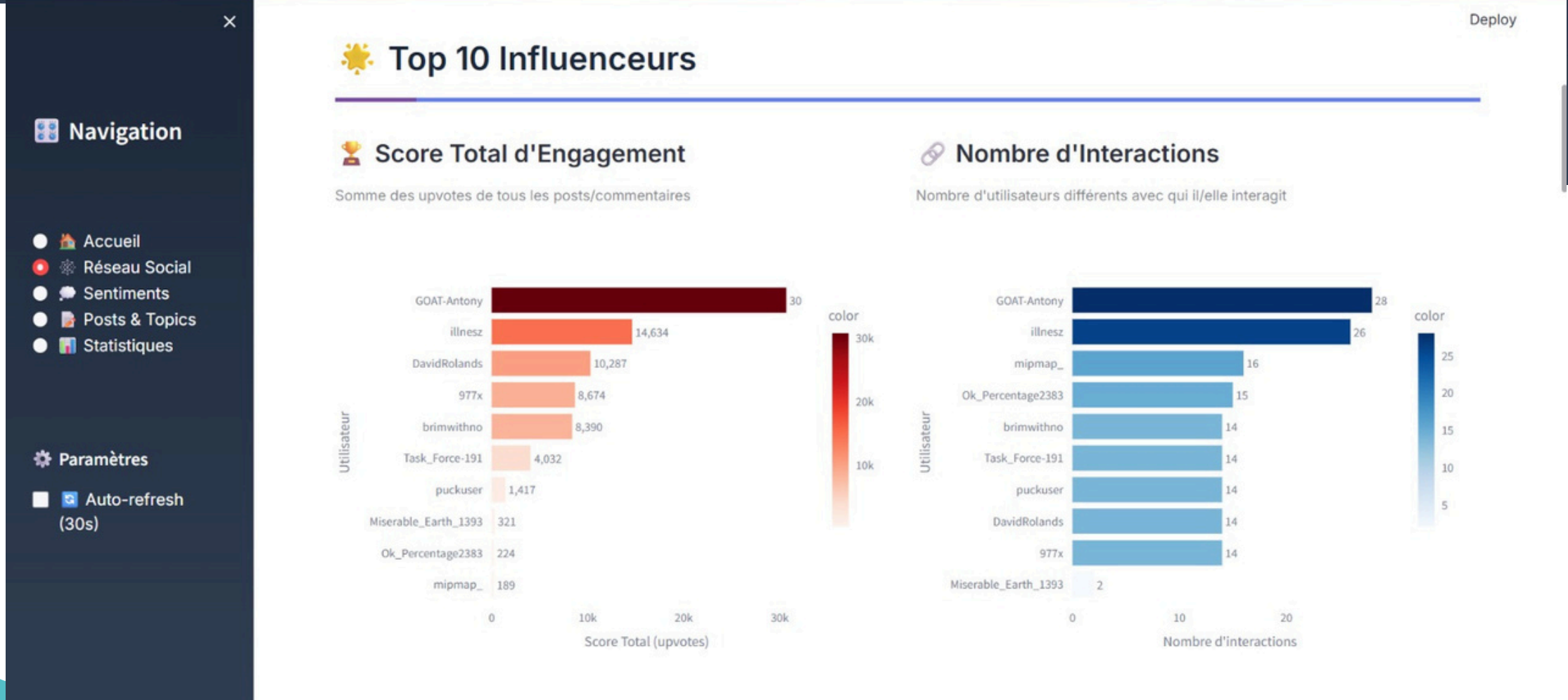
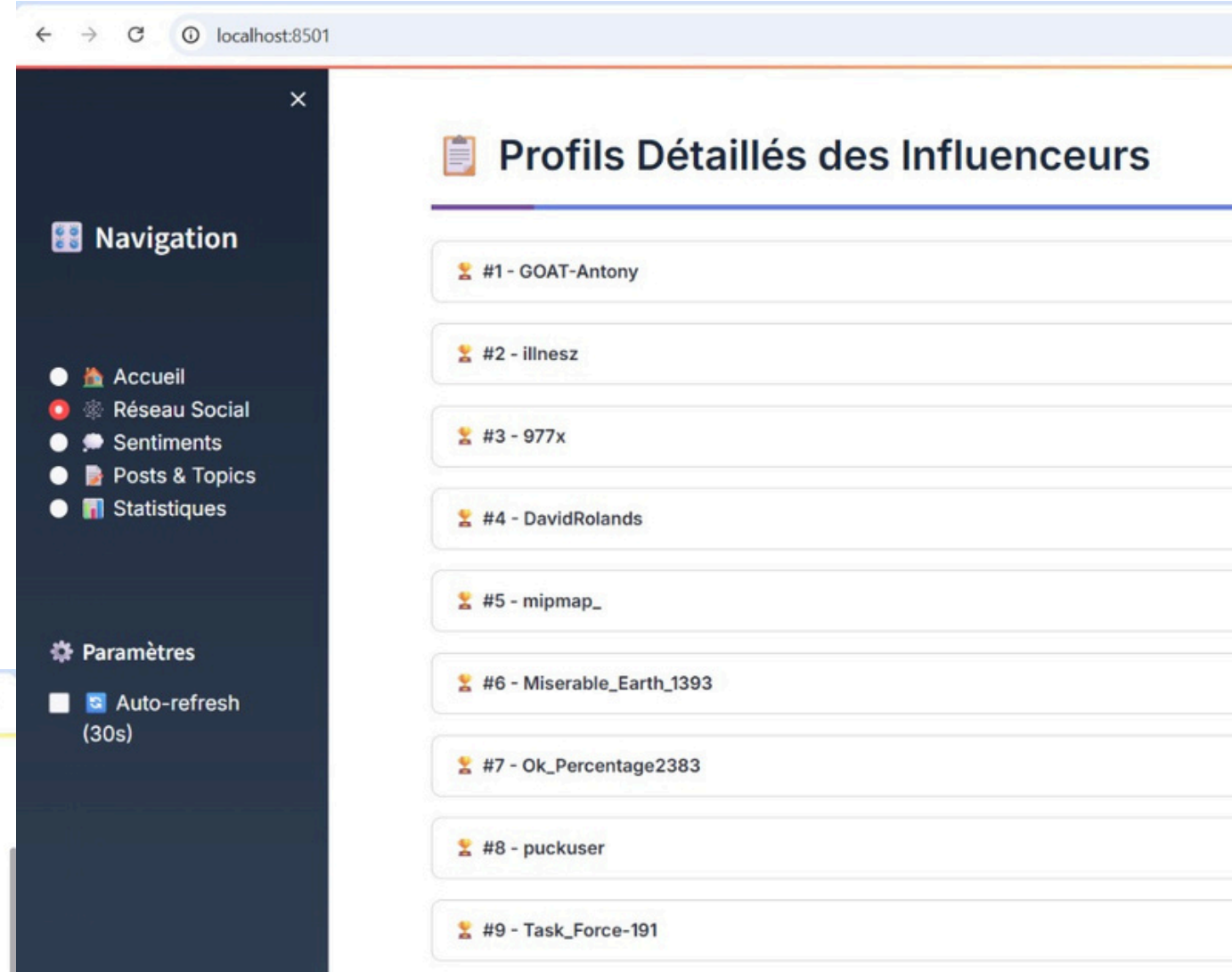


# 6 Résultats et démonstration











🏠 Communautés Détectées

Algorithme de Louvain : Détecte automatiquement les groupes d'utilisateurs qui interagissent fréquemment ensemble.

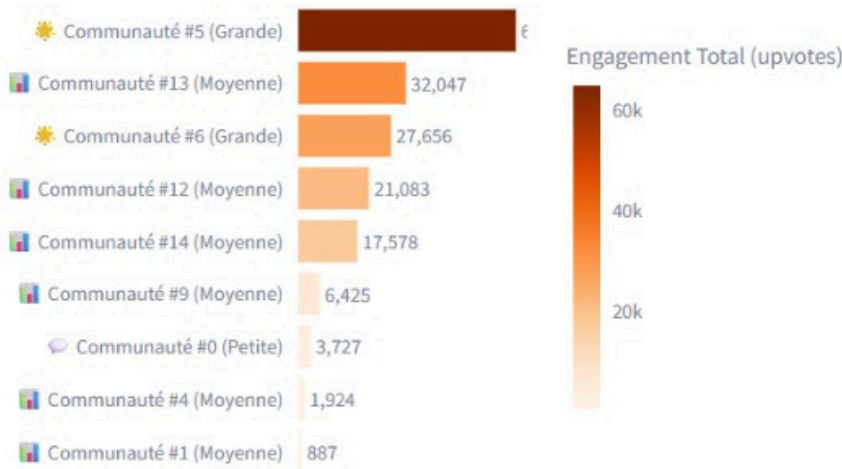
💡 Interprétation :

- Grandes communautés (>25 membres) = Groupes de fans très actifs
- Moyennes (15-25 membres) = Discussions thématiques
- Petites (<15 membres) = Niches spécialisées

👥 Taille des Communautés (Top 10)



📊 Engagement Total par Communauté



📋 Tableau Récapitulatif des Communautés

Communauté	Membres	Interactions Moy.	Engagement Total
🌟 Communauté #5 (Grande)	28	1.96	64876 🏆
🌟 Communauté #6 (Grande)	26	1.96	27656 🏆
👥 Communauté #1 (Moyenne)	17	1.88	887 🏆
👥 Communauté #2 (Moyenne)	16	1.88	598 🏆
👥 Communauté #4 (Moyenne)	15	1.87	1924 🏆
👥 Communauté #9 (Moyenne)	15	1.87	6425 🏆
👥 Communauté #12 (Moyenne)	15	1.87	21083 🏆
👥 Communauté #13 (Moyenne)	15	1.93	32047 🏆
👥 Communauté #14 (Moyenne)	15	1.93	17578 🏆
💡 Communauté #0 (Petite)	14	1.86	3727 🏆

# 🧠 Analyse des Sentiments

Classification ML des opinions

Modèle

**Random F...**

Accuracy

**85-95%**

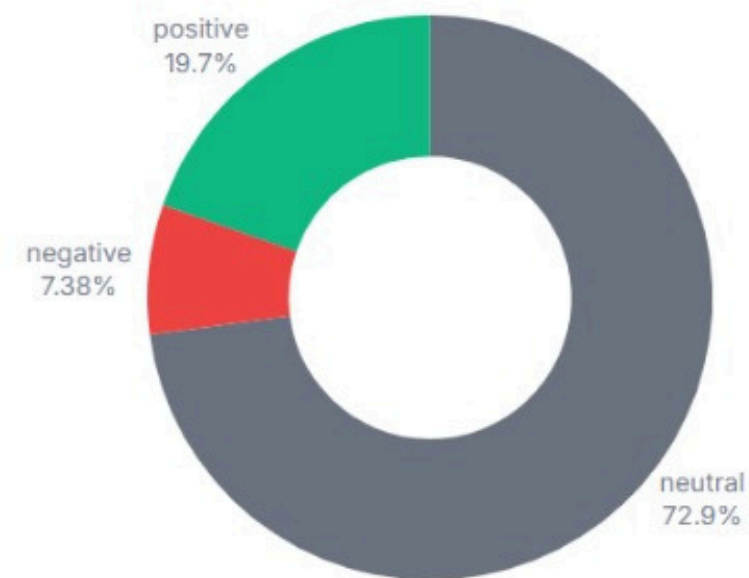
Baseline

**VADER**

Coverage

**38.7%**

## 📊 Distribution



■ neutral  
■ positive  
■ negative



Exemples

Positifs Neutres Négatifs

✓

Malian fans trolling a Tunisian fan following Mali's win vs Tunisia th...

▼

✓

Ismail Jakobs assists Nicolas Jackson for Senegal's first goal in AFCON...

▼

✓

Today's MotD: Uganda vs Nigeria (AFCON, Morocco)...

▲

★18

💜5

📌r/soccer

\*\*Oya! How far?\*\*

Today's Match of the Day is a series covering a fixture happening today, somewhere in the world, between historical rivals, across

The point is simple: you don't need a bet to care about a football match; All you really need is context.

You can follow the game on [QFAX](https://qfax.football) or your preferred scores app (Fotmo

Exemples

Positifs Neutres Négatifs

⚠️ [BBC] 'I can kill you right now' - Sudan's footballers on civil war...

▲

★2054

💜92

📌r/soccer

"They didn't even give him a chance. They shot him more than 20 or 25 times.

"One of our childhood friends was also with them, but he couldn't say anything. So he just saw our friend die in front of his eyes

The matter of fact way in which Sudan forward John Mano recounts the death of his best friend Medo is at odds with the intense loo

Med

⚠️ After poor results at Afcon (africa), the Gabon foopball team has been...

▼

⚠️ Sudan's warriors fall fighting as teenage Mbaye seals Senegal's place ...

▼

⚠️ Sudan's warriors fall fighting as teenage Mbaye seals Senegal's place

▼





# Posts & Topics

Exploration des discussions

Subreddit

Tous

Score Min

0

11530

Trier par

Score

Score

Date

Commentaires

844 posts trouvés

## Posts Détaillés

- A DR Congo fan stood for 90+ minutes during the AFCON game vs Senegal, honoring ...
- Zinedine Zidane at AFCON watching his son Luca Zidane play for Algeria...
- A fan in the Senegal vs Botswana game (no idea for which team)...



# Statistiques Avancées

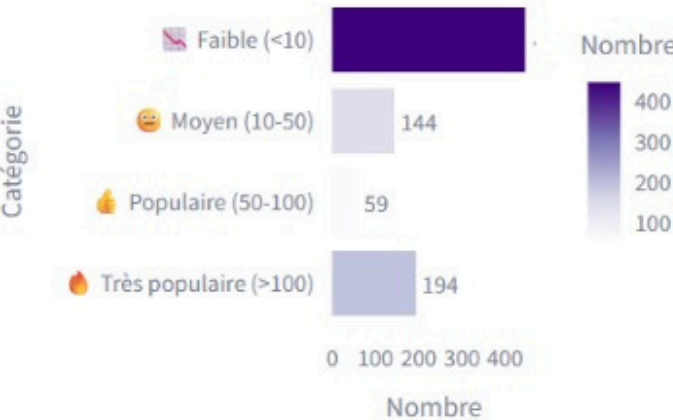
Analyses statistiques approfondies avec insights visuels



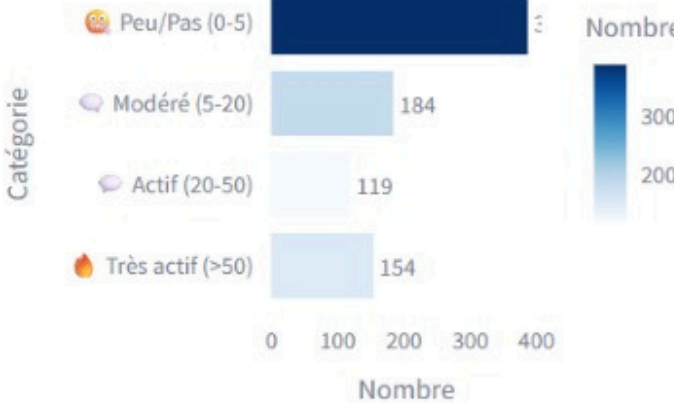
## Analyse des Données Clés



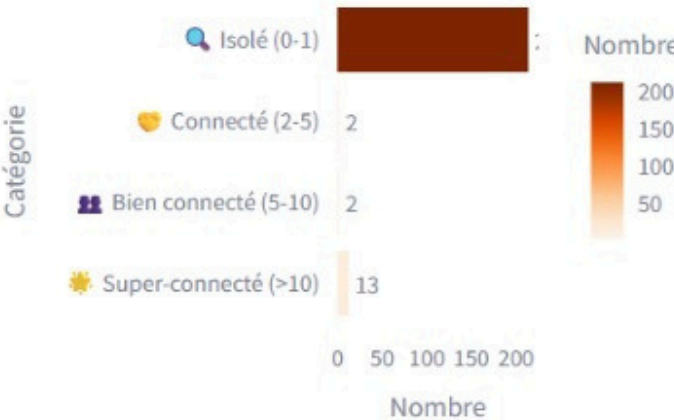
### Répartition des Scores



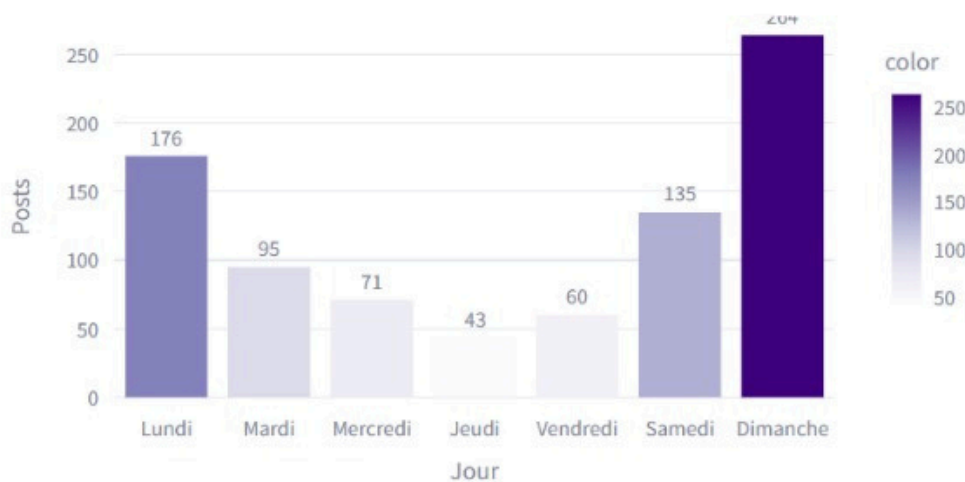
### Répartition des Commentaires



### Répartition des Connexions

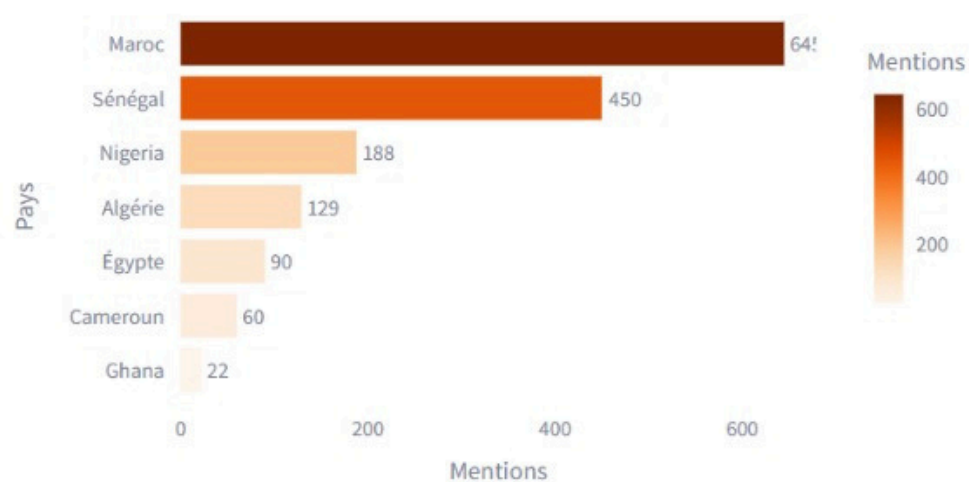


## Activité par Jour de la Semaine



🏆 Dimanche est le jour le plus actif !

## Pays Africains Mentionnés



🏆 Maroc domine avec 645 mentions !

## 🎯 Insights Statistiques

### 📊 Répartition :

- 🔥 **Top 25%** (score > 80) : 211 posts (25.0%)
- 📉 **Bottom 25%** (score < 1) : 37 posts (4.4%)
- 💡 Les posts très populaires (>100 upvotes) sont **rare**s mais génèrent beaucoup d'attention

### 💬 Engagement :

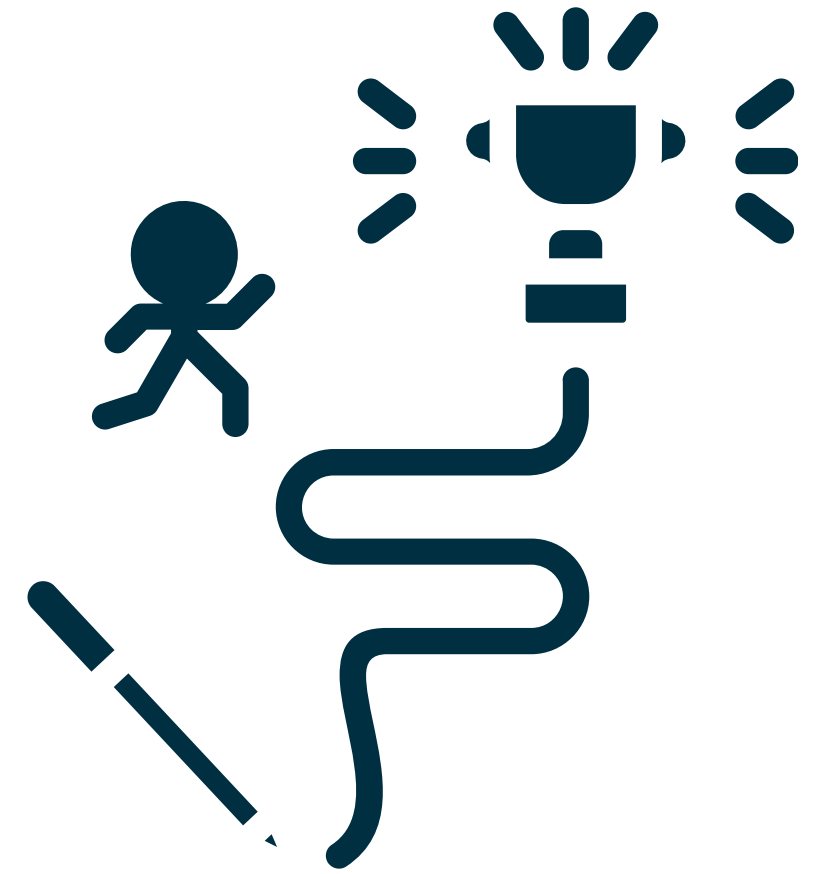
- 🔥 **Top 10%** (>81 comments) : 84 posts
- 💬 Ces posts génèrent **28,933** commentaires (72.4% du total)
- 💡 **10.0% des posts** → **72.4% des discussions**

### ☀️ Connectivité :

- ☀️ **Top 10%** (>1 connexions) : 17 utilisateurs
- 💡 Ces **super-connecteurs** sont les influenceurs du réseau
- 📊 Le reste (214 users) a <1 connexions

7

# DIFFICULTÉS RENCONTRÉES & SOLUTIONS



## RATE LIMITING REDDIT

### **Problème :**

Blocage de l'API Reddit

Risque de blocage du scraper

### **Solution :**

- Scraping espacé (30 secondes)
- Rotation User-Agent
- Gestion erreurs 429 (Too Many Req.) |

## VOLUME INITIAL LIMITÉ

### **Problème :**

Peu de posts en début de projet

Dataset insuffisant pour ML

### **Solution :**

- Élargissement des mots-clés
- Collecte historique (7 jours)
- Extraction des commentaires
- Multi-subreddits

## LABELLISATION ML

### **Problème :**

Pas de labels humains disponibles

Annotation manuelle trop coûteuse

### **Solution :**

- VADER pour génération auto labels
- Validation sur échantillon manuel
- Approche semi-supervisée

7

# Conclusion et perspective





## Objectifs atteints



- Architecture Kappa fonctionnelle et performante
- Pipeline temps réel avec latence 1–2 minutes
- ML fiable (85–95% accuracy)
- Réseau social cartographié (15 communautés, 20 influenceurs)
- Dashboard interactif opérationnel
- Résultats concrets sur données réelles
- Projet applicable à des cas métiers réels

## PERSPECTIVES D'ÉVOLUTION



Topic Modeling (LDA/BERTopic)

Alertes Temps Réel

API Reddit Officielle

Deep Learning (BERT/RoBERTa)

Détection Anomalies

**Merci pour votre  
attention**

