

Projet NoSQL : Kafka + Elasticsearch + Kibana

Un pipeline de données véhicules connectés

Amine, Samy, Mehdi
Étudiants BDML, EFREI

3 avril 2025

Table des matières

1	Introduction	2
2	Installation	4
2.1	Prérequis	4
2.2	Étapes d'installation et de configuration	4
2.2.1	Captures d'écran (Installation)	5
3	Création et manipulation des données (CRUD)	6
3.1	Mise en place du producteur de données	6
3.2	Consommation et insertion dans Elasticsearch	7
3.3	CRUD : Créer, Lire, Mettre à jour, Supprimer	8
3.3.1	Créer (Create)	8
3.3.2	Lire (Read)	8
3.3.3	Mettre à jour (Update)	8
3.3.4	Supprimer (Delete)	8
4	Fonctionnalités avancées : agrégations et visualisations	9
4.1	Exemple de dashboard	9
5	Conclusion et recommandations	10
5.1	Avantages de notre pipeline	10
5.2	Limites	10
5.3	Comparaison succincte	10
5.4	Perspectives	10
5.5	Mot de la fin	11

Chapitre 1

Introduction

Dans le cadre de notre projet NoSQL, nous avons choisi de travailler sur un écosystème combinant plusieurs technologies :

- **Apache Kafka** : une plateforme de streaming distribuée, open source, initialement développée par LinkedIn puis gérée par la fondation Apache. Elle est principalement écrite en Java et Scala.
- **Elasticsearch** : un moteur de recherche et d'analytics temps réel très performant, développé en Java par la société Elastic. Il est historiquement open source (sous licence Apache), puis a évolué vers la *Elastic License 2.0* (certaines fonctionnalités payantes).
- **Kibana** : un outil de visualisation et d'exploration des données pour Elasticsearch, également proposé par Elastic.

Notre projet consiste à mettre en place un **pipeline de données** où nous simulons des véhicules (positions GPS, vitesse, niveau de batterie, etc.) afin de tester la scalabilité, la résilience et les fonctionnalités de ces technologies. Nous avons choisi ces outils car ils sont très répandus dans le monde de la data (ingestion temps réel, indexation, analytics, monitoring...) et ils couvrent un large panel de cas d'usage : ingestion massive, stockage distribué, recherche full-text, création de dashboards, etc.

Comparaison avec d'autres NoSQL (MongoDB, Cassandra, Neo4j, Redis) :

Technologie	Modèle	Points forts	Points faibles
MongoDB	Document (JSON)	<ul style="list-style-type: none">— Très flexible— Facile à installer	<ul style="list-style-type: none">— Moins efficace pour la recherche full-text (comparé à Elastic)— Peut nécessiter un schéma plus rigoureux pour les grosses apps

Cassandra	Colonnes distribuées	<ul style="list-style-type: none"> — Très scalable en écriture — Faible latence 	<ul style="list-style-type: none"> — Langage de requête (CQL) plus limité que SQL — Pas idéal pour la recherche textuelle
Neo4j	Graphes	<ul style="list-style-type: none"> — Excellente gestion des relations — Cypher = requêtes graphiques puissantes 	<ul style="list-style-type: none"> — Moins adapté à l'ingestion massive temps réel — Moins efficace pour l'analytics basique
Redis	Clé-Valeur en mémoire	<ul style="list-style-type: none"> — Extrêmement rapide — Bonne pour le cache 	<ul style="list-style-type: none"> — Stockage en RAM = coûts élevés — Pas d'analytics avancé
Kafka + Elasticsearch	<i>(Middleware + Moteur de recherche)</i>	<ul style="list-style-type: none"> — Ingestion temps réel (Kafka) — Recherche/analytics (Elastic) — Visualisation (Kibana) 	<ul style="list-style-type: none"> — Plus complexe à mettre en place (plusieurs services) — Infrastructure distribuée (maintenance, monitoring)

Ainsi, Kafka + Elasticsearch + Kibana présentent un excellent compromis entre ingestion temps réel et analyses rapides, notamment sur des volumes importants de données.

Chapitre 2

Installation

2.1 Prérequis

- Docker Desktop (ou Docker Engine) installé sur la machine.
- Python 3 (pour simuler les données).

2.2 Étapes d'installation et de configuration

Nous avons créé un fichier `docker-compose.yml` permettant de lancer :

- **Zookeeper** (nécessaire pour Kafka)
- **Kafka**
- **Elasticsearch**
- **Kibana**

Voici le contenu principal (résumé) :

Listing 2.1 – Extrait du `docker-compose.yml`

```
version: '3.8'
services:
  zookeeper:
    image: confluentinc/cp-zookeeper:7.5.0
    ...

  kafka:
    image: confluentinc/cp-kafka:7.5.0
    ...

  elasticsearch:
    image: docker.elastic.co/elasticsearch/elasticsearch:8.11.1
    environment:
      - discovery.type=single-node
      - xpack.security.enabled=false
    ports:
      - "9200:9200"
    networks:
      - elk

  kibana:
    image: docker.elastic.co/kibana/kibana:8.11.1
    environment:
```

```

    - ELASTICSEARCH_HOSTS=http://elasticsearch:9200
ports:
    - "5601:5601"
networks:
    - elk

```

networks:
elk:

Une fois ce fichier en place, on lance la commande : `docker-compose up -d`

2.2.1 Captures d'écran (Installation)

Nous avons réalisé des captures d'écran pendant l'installation. Ci-dessous, **nous incluons les images sans réécrire leur contenu**, conformément aux consignes :

```

• (base) aminemzali@MacBook-Pro-de-Sylvain ProjetNoSQL % mkdir kafka-vehicle-project
cd kafka-vehicle-project
touch docker-compose.yml
• (base) aminemzali@MacBook-Pro-de-Sylvain kafka-vehicle-project % docker-compose up -d
WARN[0000] /Users/aminemzali/Documents/Efrei/S8/ProjetNoSQL/kafka-vehicle-project/docker-compose.yml: the attribute `version` is obsolete, it will
be ignored, please remove it to avoid potential confusion
[+] Running 39/39
  ✓ elasticsearch Pulled          91.2s
  ✓ zookeeper Pulled            101.9s
  ✓ kibana Pulled               69.3s
  ✓ kafka Pulled               101.9s
[+] Running 6/6
  ✓ Network kafka-vehicle-project_elk Created      0.1s
  ✓ Network kafka-vehicle-project_default Created  0.1s
  ✓ Container zookeeper Started                   0.4s
  ✓ Container elasticsearch Started                0.4s
  ✓ Container kibana Started                      0.4s
  ✓ Container kafka Started                      0.4s
○ (base) aminemzali@MacBook-Pro-de-Sylvain kafka-vehicle-project %

```

FIGURE 2.1 – Création du dossier `kafka-vehicle-project` et exécution de `docker-compose up -d`

```

• (base) aminemzali@MacBook-Pro-de-Sylvain kafka-vehicle-project % mkdir producer
cd producer
touch simulate_vehicles.py
• (base) aminemzali@MacBook-Pro-de-Sylvain producer % python3 -m venv venv
source venv/bin/activate
pip install kafka-python faker
Collecting kafka-python
  Downloading kafka_python-2.1.4-py2.py3-none-any.whl (276 kB)
    ━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 276.6/276.6 kB 4.1 MB/s eta 0:00:00
Collecting faker
  Downloading faker-37.1.0-py3-none-any.whl (1.9 MB)
    ━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 1.9/1.9 MB 12.0 MB/s eta 0:00:00
Collecting tzdata
  Using cached tzdata-2025.2-py2.py3-none-any.whl (347 kB)
Installing collected packages: kafka-python, tzdata, faker
Successfully installed faker-37.1.0 kafka-python-2.1.4 tzdata-2025.2

[notice] A new release of pip is available: 23.0.1 -> 25.0.1
[notice] To update, run: pip install --upgrade pip
○ (venv) (base) aminemzali@MacBook-Pro-de-Sylvain producer %

```

FIGURE 2.2 – Installation de Python `kafka-python` et `faker`

Chapitre 3

Création et manipulation des données (CRUD)

3.1 Mise en place du producteur de données

Pour simuler des données de véhicules, nous avons un script `simulate_vehicles.py` qui envoie des messages vers Kafka.

— **Création du topic :**

```
docker exec -it kafka \
  kafka-topics --create --topic vehicle-data \
  --bootstrap-server localhost:9092 \
  --partitions 1 \
  --replication-factor 1
```

— **Producteur Python :**

Listing 3.1 – Extrait de `simulate_vehicles.py`

```
from kafka import KafkaProducer
from faker import Faker
import json, time, random

fake = Faker()
producer = KafkaProducer(
    bootstrap_servers='localhost:9092',
    value_serializer=lambda v: json.dumps(v).encode('utf-8')
)

def generate_vehicle_data(vehicle_id):
    return {
        "vehicle_id": vehicle_id,
        "@timestamp": time.strftime("%Y-%m-%dT%H:%M:%SZ", time.gmtime()),
        "latitude": float(fake.latitude()),
        "longitude": float(fake.longitude()),
        "speed_kmh": round(random.uniform(0, 120), 2),
        "battery_level": random.randint(10, 100)
    }

if __name__ == "__main__":
    vehicle_ids = [f"veh-{i}" for i in range(1, 6)]
```

```

while True:
    vehicle = random.choice(vehicle_ids)
    data = generate_vehicle_data(vehicle)
    print(f"Sending: {data}")
    producer.send("vehicle-data", data)
    time.sleep(1)

```

Quelques extraits de la console :

```

el': 57}
Sending: {'vehicle_id': 'veh-4', 'timestamp': 1743664644.611537, 'latitude': 34.305018, 'longitude': 143.505552, 'speed_kmh': 59.36, 'battery_level': 99}
Sending: {'vehicle_id': 'veh-2', 'timestamp': 1743664645.61519, 'latitude': -25.614558, 'longitude': 89.926027, 'speed_kmh': 39.91, 'battery_level': 93}
Sending: {'vehicle_id': 'veh-4', 'timestamp': 1743664646.619885, 'latitude': 6.4787735, 'longitude': 177.222666, 'speed_kmh': 3.92, 'battery_level': 79}
Sending: {'vehicle_id': 'veh-4', 'timestamp': 1743664647.6234848, 'latitude': -53.427157, 'longitude': 73.846512, 'speed_kmh': 104.39, 'battery_level': 43}
Sending: {'vehicle_id': 'veh-4', 'timestamp': 1743664648.6288981, 'latitude': 39.9791005, 'longitude': -22.054631, 'speed_kmh': 15.1, 'battery_level': 75}
Sending: {'vehicle_id': 'veh-5', 'timestamp': 1743664649.631018, 'latitude': 49.169727, 'longitude': 28.261668, 'speed_kmh': 85.02, 'battery_level': 32}
Sending: {'vehicle_id': 'veh-2', 'timestamp': 1743664650.6359448, 'latitude': -72.305269, 'longitude': 159.576207, 'speed_kmh': 20.96, 'battery_level': 87}
Sending: {'vehicle_id': 'veh-5', 'timestamp': 1743664651.6398091, 'latitude': -41.6065435, 'longitude': -2.282485, 'speed_kmh': 47.29, 'battery_level': 57}

```

FIGURE 3.1 – Console du producteur : envoi des données simulées

3.2 Consommation et insertion dans Elasticsearch

Pour lire les données depuis Kafka et les envoyer à Elasticsearch, nous utilisons un `consumer_to_es.py` :

Listing 3.2 – Extrait de `consumer_to_es.py`

```

from kafka import KafkaConsumer
from elasticsearch import Elasticsearch
import json

# Connexion Elasticsearch
es = Elasticsearch("http://localhost:9200")

# Connexion au topic Kafka
consumer = KafkaConsumer(
    "vehicle-data",
    bootstrap_servers="localhost:9092",
    value_deserializer=lambda v: json.loads(v.decode("utf-8")),
    auto_offset_reset='earliest',
    enable_auto_commit=True,
    group_id="vehicle-consumer-group"
)

for message in consumer:
    data = message.value
    print(f"Received: {data}")
    # Indexation
    es.index(index="vehicle-data", document=data)

```



```
(venv) (base) aminemzali@MacBook-Pro-de-Sylvain producer % python consumer_to_es.py
Received: {'vehicle_id': 'veh-1', '@timestamp': '2025-04-03T07:31:50Z', 'latitude': -44.7262915, 'longitude': -1
58.483988, 'speed_kmh': 3.59, 'battery_level': 13}
Received: {'vehicle_id': 'veh-2', '@timestamp': '2025-04-03T07:31:51Z', 'latitude': -47.789725, 'longitude': 116
.624331, 'speed_kmh': 20.68, 'battery_level': 50}
Received: {'vehicle_id': 'veh-4', '@timestamp': '2025-04-03T07:31:52Z', 'latitude': 58.216725, 'longitude': -142
.907469, 'speed_kmh': 35.53, 'battery_level': 10}
Received: {'vehicle_id': 'veh-4', '@timestamp': '2025-04-03T07:31:53Z', 'latitude': 59.4942935, 'longitude': 83.
211519, 'speed_kmh': 17.69, 'battery_level': 69}
Received: {'vehicle_id': 'veh-4', '@timestamp': '2025-04-03T07:31:54Z', 'latitude': 75.139365, 'longitude': 64.2
12467, 'speed_kmh': 39.71, 'battery_level': 75}
Received: {'vehicle_id': 'veh-5', '@timestamp': '2025-04-03T07:31:55Z', 'latitude': -15.007959, 'longitude': 78.
296096, 'speed_kmh': 94.04, 'battery_level': 86}
Received: {'vehicle_id': 'veh-1', '@timestamp': '2025-04-03T07:31:56Z', 'latitude': 31.4283075, 'longitude': 35.
395944, 'speed_kmh': 73.6, 'battery_level': 10}
Received: {'vehicle_id': 'veh-4', '@timestamp': '2025-04-03T07:31:57Z', 'latitude': -81.028558, 'longitude': 16.
888941, 'speed_kmh': 86.07, 'battery_level': 96}
Received: {'vehicle_id': 'veh-5', '@timestamp': '2025-04-03T07:31:58Z', 'latitude': -74.770449, 'longitude': -40
.474780, 'speed_kmh': 96.00, 'battery_level': 60}
```

FIGURE 3.2 – Console du consommateur : réception et indexation dans Elasticsearch

3.3 CRUD : Créer, Lire, Mettre à jour, Supprimer

3.3.1 Créer (Create)

L'instruction `es.index(index="vehicle-data", document=data)` permet de *créer* automatiquement un document dans l'index `vehicle-data`. Si l'index n'existe pas, il est créé dynamiquement.

3.3.2 Lire (Read)

Nous pouvons lire les données via Kibana (section “Discover”) ou via une requête REST :

```
curl -X GET http://localhost:9200/vehicle-data/_search?pretty
```

3.3.3 Mettre à jour (Update)

Pour modifier un document précis (connaissant son `_id`) :

```
curl -X POST "localhost:9200/vehicle-data/_update/<document_id>" \
-H 'Content-Type: application/json' -d '{
  "doc": {
    "speed_kmh": 42.0
  }
}'
```

3.3.4 Supprimer (Delete)

```
curl -X DELETE "localhost:9200/vehicle-data/_doc/<document_id>"
```

Ou pour supprimer l'index complet :

```
curl -X DELETE http://localhost:9200/vehicle-data
```

Chapitre 4

Fonctionnalités avancées : agrégations et visualisations

Nous utilisons Kibana pour créer des **Index Patterns** (Data Views) et des **dashboards**. Une fois que **vehicle-data** est créé et que le champ **@timestamp** est reconnu comme date, nous pouvons :

- Explorer les données dans Kibana (Discover)
- Créer des visualisations (histogramme, pie chart, bar chart, etc.)

4.1 Exemple de dashboard

Nous avons créé trois visualisations :

1. Un **histogramme** du nombre de messages par tranche de temps
2. Un **pie chart** pour la répartition par **vehicle_id**
3. Un **bar chart** de la vitesse moyenne par véhicule

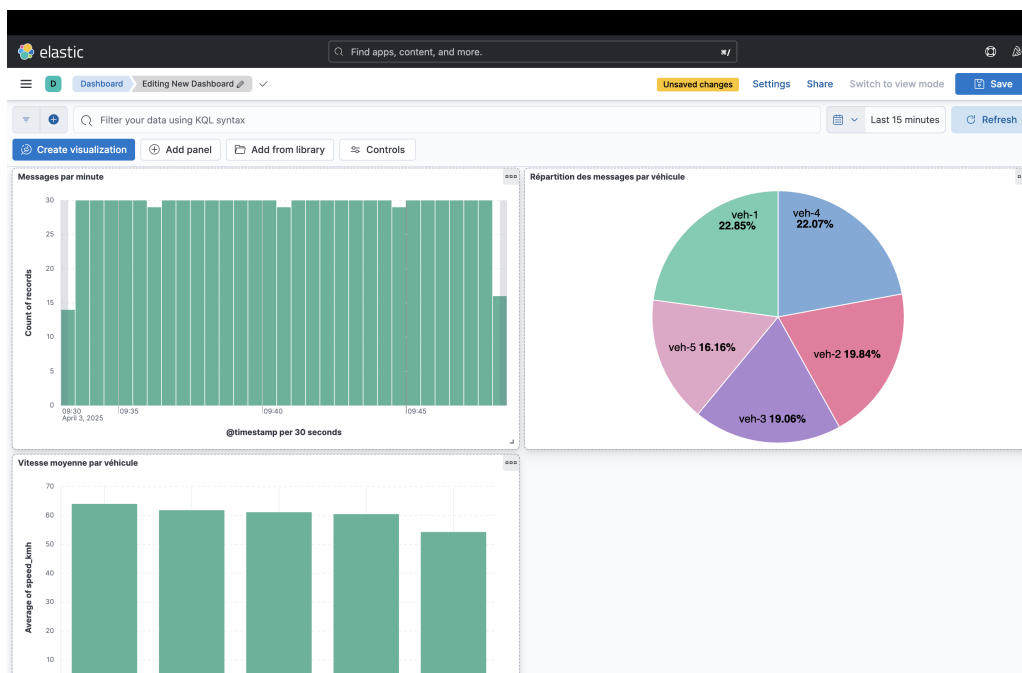


FIGURE 4.1 – Dashboard dans Kibana : histogramme, pie chart et bar chart

Chapitre 5

Conclusion et recommandations

5.1 Avantages de notre pipeline

- **Ingestion temps réel** : Kafka gère un flux continu de données et permet de scaler horizontalement en augmentant le nombre de partitions.
- **Indexation et recherche rapide** : Elasticsearch permet de chercher dans les documents quasi instantanément, et de faire des agrégations complexes.
- **Visualisation ergonomique** : Kibana fournit une interface simple pour créer des dashboards, analyser les données et effectuer des recherches.

5.2 Limites

- L'**infrastructure** est plus lourde que des solutions “simples” (ex. MongoDB). Il faut maintenir Kafka, Elasticsearch, Kibana, Zookeeper.
- Des **coûts d'hébergement** plus élevés dès lors qu'on veut de la haute disponibilité (clusters distribués).

5.3 Comparaison succincte

En comparaison à MongoDB, Cassandra, Neo4j ou Redis, ce pipeline met l'accent sur :

- **Consistance et scalabilité** (Kafka + ES sont distribués)
- **Recherche avancée et analytics temps réel** (Elasticsearch excelle en full-text et agrégations)
- **Résistance aux pannes** via la réplication Kafka et Elasticsearch.

D'autres solutions NoSQL (MongoDB, Redis, etc.) sont parfois plus simples à installer et maintenir, mais moins performantes pour un cas “ingestion massive + requêtes complexes + visualisation temps réel”.

5.4 Perspectives

Pour aller plus loin, nous pourrions :

- Ajouter un **Kafka Connect** pour automatiser l'ingestion vers Elasticsearch
- Jouer sur le **nombre de partitions** Kafka et de **shards/réplicas** ES pour mesurer les performances
- Gérer la **géolocalisation** avec un type `geo_point` pour afficher les véhicules sur une carte.

5.5 Mot de la fin

Ce pipeline répond aux besoins d'une entreprise souhaitant *collecter et analyser en temps réel* une grande quantité de données de capteurs ou de logs. Grâce à l'élasticité de Kafka et la puissance d'Elasticsearch, associées à la visualisation de Kibana, nous obtenons une solution performante, évolutive et relativement simple à déployer via Docker.

— **Fin du Rapport** —