
PROJET FINAL NFE 101 2024 - 2025

Réalisé par : FONKOUO SANOU TEMENA LOIC

Introduction

Dans le cadre de l'Unité d'Enseignement **NFE 101**, portant sur l'intégration et l'évolution des systèmes d'information, nous avons été amenés à concevoir et développer un projet personnel basé sur les exemples étudiés en cours. Ce projet vise à mettre en pratique les concepts clés abordés, en explorant les capacités d'intégration et de traitement des flux de données en temps réel à l'aide de **Kafka**.

Pour répondre aux exigences du projet, nous avons développé plusieurs mini-applications Kafka, chacune jouant un rôle spécifique dans le traitement, la transformation ou l'analyse des données.

Afin de mieux situer le contexte, cette introduction s'articulera autour des éléments suivants :

1. **Description du jeu de données utilisé** : Nous présenterons en détail les données traitées, leur origine, ainsi que leur structure.
2. **Travail de prétraitement des données** : Nous expliquerons les étapes et les techniques employées pour préparer les données avant leur ingestion dans Kafka.
3. **Architecture de notre application** : Nous présenterons à l'aide d'un schéma l'architecture de nos applications, pour illustrer leur rôle.
4. **Description des mini-applications Kafka** : Nous décrirons le rôle et le fonctionnement de chaque application, en les illustrant par des captures d'écran et des explications techniques pertinentes.

Description du jeu de données

Le jeu de données intitulé "[Carte des pharmacies de Paris](#)" fournit une liste détaillée des pharmacies situées à Paris, incluant leurs coordonnées, adresses et informations de géolocalisation. Ces données proviennent du répertoire FINESS, avec une mise à jour datant de septembre 2013.

Le fichier contient 987 enregistrements, chacun représentant une pharmacie, avec les champs suivants :

Champ	Description
nofinesset	Numéro FINESS de l'établissement.
nofinesselj	Numéro FINESS de l'entité juridique.
rs	Raison sociale de la pharmacie.
rslongue	Raison sociale longue.
complidistrib	Complément de distribution.
numvoie	Numéro de la voie.
typvoie	Type de voie (rue, avenue, boulevard, etc.).
voie	Libellé de la voie.
compvoie	Complément de voie.
lieuditbp	Lieu-dit ou Boîte Postale.
departement	Code du département.
libdepartement	Libellé du département.
cp	Code postal.
commune	Nom de la commune.
telephone	Numéro de téléphone.
telecopie	Numéro de télécopie.
dateouv	Date d'ouverture de la pharmacie.
dateautor	Date d'autorisation.
datemaj	Date de la dernière mise à jour des informations.
wgs84	Coordonnées géographiques au format WGS84 (latitude, longitude).
lat	Latitude.
lng	Longitude.

Figure 1: Description des colonnes du jeu de données.

Ce jeu de données est disponible aux formats CSV et JSON, facilitant ainsi son intégration et son exploitation dans divers systèmes d'information. Il est distribué sous la **Licence Ouverte / Open Licence version 2.0**, permettant une utilisation libre sous réserve de mentionner la source.

Travail de prétraitement des données

Pour garantir la qualité des données et leur intégration dans notre application, nous avons effectué un prétraitement minutieux en utilisant plusieurs outils adaptés. Voici les étapes réalisées :

- Correction des espaces et des colonnes mal renseignées

Nous avons utilisé OpenRefine pour détecter et corriger les espaces mal renseignés dans certaines colonnes, notamment celles contenant des données de type date. Cela a permis d'homogénéiser les formats de ces colonnes et de garantir une meilleure cohérence.

- 2. Gestion des caractères spéciaux liés à l'encodage

Certains caractères spéciaux provenant d'un problème d'encodage dans les données brutes ont été identifiés. OpenRefine a également été utilisé pour les corriger, en remplaçant ces caractères par des alternatives correctes.

- 3. Problèmes de décalage dans les colonnes

Le jeu de données contenait certaines lignes problématiques avec des virgules en surplus, ce qui entraînait un désalignement entre les colonnes et leurs en-têtes. Pour résoudre ce problème, nous avons utilisé Visual Studio Code pour rechercher rapidement les lignes contenant un nombre d'occurrences de virgules supérieur à celui attendu. Cela nous a permis de corriger ces lignes manuellement.

- Changement du séparateur de colonnes

Initialement, le fichier CSV utilisait le point-virgule (;) comme séparateur, ce qui a posé des problèmes lors de l'importation des données dans notre classe CsvReader. Nous avons tenté de changer le séparateur de

point-virgule à virgule (,), à l'aide d'Excel, mais cela n'a pas permis de résoudre le problème. Finalement, nous avons conservé le point-virgule comme séparateur et avons ajusté notre code en précisant explicitement ce séparateur lors de la lecture des données.

En résumé, bien que le jeu de données ne présentât pas de problèmes majeurs nécessitant un prétraitement exhaustif, ces étapes ont permis de garantir la qualité et la compatibilité des données pour leur exploitation dans notre application.

Architecture de notre application

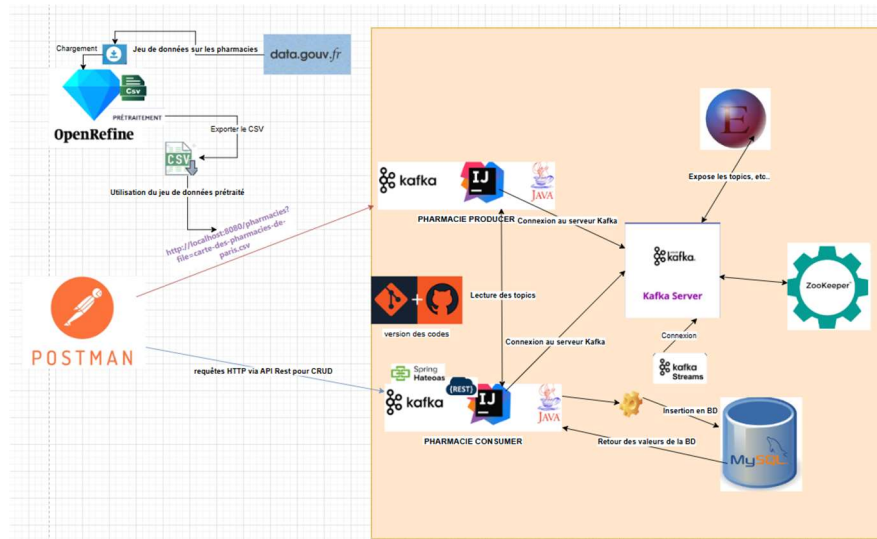


Figure 2: Architecture de notre application.

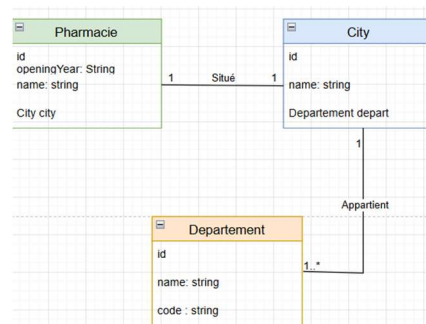


Figure 3: Diagramme de classe.

1. Obtention des Données

Les données brutes ont été obtenues à partir des topics Kafka où elles sont publiées.

2. Prétraitement et Structuration

Une fois les données récupérées, un processus de prétraitement est appliqué pour corriger, valider et structurer les informations (comme mentionné dans la section précédente). Cela garantit que seules des données propres et conformes aux exigences du système sont intégrées.

3. Mapping Objet-Relationnel (ORM)

Après le prétraitement, les données sont transformées en objets métier grâce à un processus de Mapping Objet-Relationnel (ORM). Ces objets sont ensuite reliés aux tables correspondantes de la base de données, suivant la structure définie dans le diagramme de classe présenté plus haut.

4. Stockage en Base de Données

Les objets mappés sont finalement persistés dans la base de données. Cette opération permet de stocker efficacement les informations pour une utilisation ultérieure dans d'autres parties de l'application (Ces données sont exposées sous forme d'API Rest par le Consumer).

Description des mini-applications Kafka

Pharmacies-Producer

Cette application joue le rôle de Producer dans notre architecture Kafka. Elle a pour mission de :

- Lire le fichier contenant les données nettoyées.
- Injecter chaque enregistrement de manière unitaire dans un topic Kafka spécifique.

Ce module constitue l'entrée principale du pipeline de traitement de données.

Pharmacies-Consumer

Cette application représente le Consumer principal et contient une grande partie de la logique métier. Ses principales responsabilités sont :

- Consommer les messages Kafka : Récupérer les enregistrements depuis le topic Kafka.
- Traitement des données : Effectuer les traitements nécessaires pour transformer les données en objets prêts à être persistés.
- Insertion dans la base de données : Stocker les données traitées dans une base MySQL en respectant le modèle relationnel défini.
- Exposition des données via API REST : Grâce à Spring HATEOAS, l'application fournit des endpoints CRUD (Create, Read, Update, Delete) pour interagir avec les données via des services RESTful.

Pharmacies-Streams

Initialement, les fonctionnalités de cette application étaient intégrées dans Pharmacies-Consumer. Cependant, pour réduire les couplages et améliorer la modularité, une application distincte a été mise en place.

- Cette application est dédiée à l'agrégation et à l'analyse des données en temps réel.
- Elle applique des opérations spécifiques (telles que le regroupement et le comptage) sur les messages provenant du topic Kafka.

Client Postman

Cet outil est utilisé pour tester et valider les endpoints REST exposés par l'application Pharmacies-Consumer. Il permet de simuler des appels aux API (GET, POST, PUT, DELETE) afin de s'assurer que les services RESTful fonctionnent correctement et répondent aux attentes.

Quelques captures d'écrans

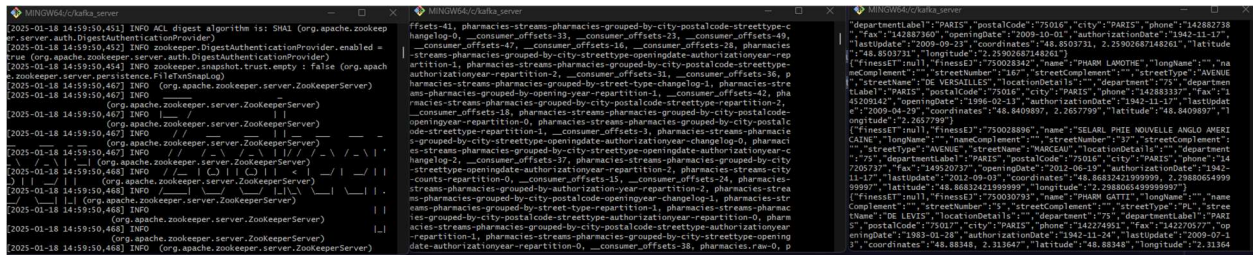


Figure 4: Server Kafka, Zookeeper, contenu des messages envoyés.

pharmacy_id	authorization_date	fax	finesse_id	finesse_id	test_update	latitude	location_details	long_name	longitude	name	name_complement	opening_date	phone	street_complement
1	1942-12-04		75008757		NULL	2019-09-29	48.8685958		2.3297288	PHARMACIE GAILOUX		1994-07-29	14388814	
2	2007-06-04		147001018	750030108	NULL	2012-04-04	48.862543	SE LAIR	2.3425080	LA PHARMACIE DE LA PLACE DE LA REPUBLIQUE		2012-01-15	147001008	
3	1943-07-06		148071538	750010530	NULL	2019-09-29	48.8605543		2.3483397	PHARMACIE PREMIERE		2009-11-16	148076230	
4	1942-11-08		145345555	750010233	NULL	2019-01-19	48.8530396		2.3537558	PHARMACIE DE LA SAINTE LUCIE		2019-01-01	145345555	
5	1942-11-11		142507275	750010485	NULL	2009-11-25	48.8391078		2.3304983	PHARMACIE DE LA MATERNITE		1991-05-30	145345555	
6	1942-11-08		142365020	750010942	NULL	2009-02-27	48.8607915		2.3480684	PHARMACIE DU VAL DE GRACE		1994-09-29	142365020	
7	1943-05-27		147073543	750014896	NULL	2009-11-21	48.8418162		2.3487131	PHARMACIE DU TEMPLE		2009-11-01	143181444	
8	1942-11-11		143177290	750014211	NULL	2019-09-29	48.8407303		2.3515551	PHARMACIE LAMON		1994-12-30	143177290	

Figure 5: Données insérées en DB.

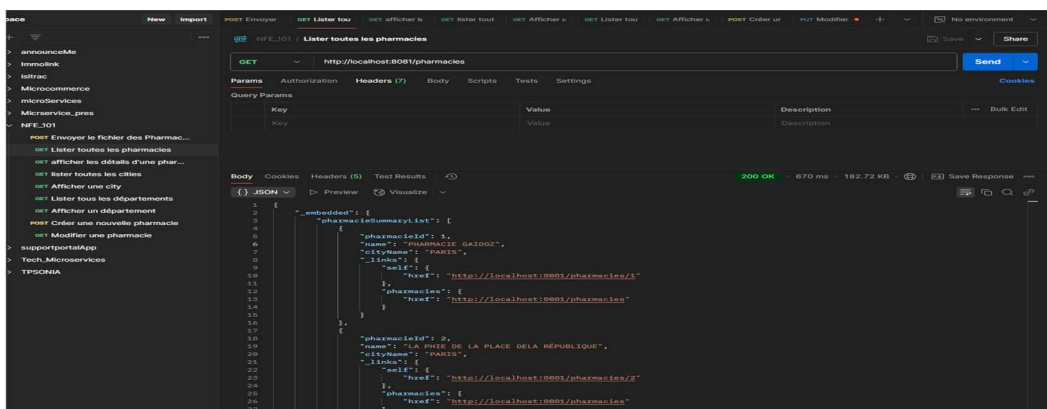


Figure 6: Liste non exhaustive des différents endpoints testés.

Timestamp	Log Level	Logger	Message	Count
2025-01-20 13:30:13	INFO	t.f.n.p.s.PharmaciesStreamProcessor	City PostalCode OpeningDate: PARIS 75008 2012-01-15	14
2025-01-20 13:30:13	INFO	t.f.n.p.s.PharmaciesStreamProcessor	City PostalCode OpeningDate: PARIS 75008 1994-09-29	14
2025-01-20 13:30:13	INFO	t.f.n.p.s.PharmaciesStreamProcessor	City PostalCode OpeningDate: PARIS 75006 1975-09-19	14
2025-01-20 13:30:13	INFO	t.f.n.p.s.PharmaciesStreamProcessor	City PostalCode OpeningDate: PARIS 75007 2012-09-03	14
2025-01-20 13:30:13	INFO	t.f.n.p.s.PharmaciesStreamProcessor	City PostalCode OpeningDate: PARIS 75009 2010-10-15	14
2025-01-20 13:30:13	INFO	t.f.n.p.s.PharmaciesStreamProcessor	City PostalCode OpeningDate: PARIS 75010 2007-04-30	14

Figure 7: Exemple de groupement par City|Codepostal|OpeningDate.


```
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - Authorization Year: 1944 Count: 210
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - Authorization Year: 1960 Count: 14
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - Authorization Year: 1977 Count: 28
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - Authorization Year: 1942 Count: 6762
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75010|PL Count: 14
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75007|PL Count: 14
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75014|PL Count: 42
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75002|BOULEVARD Count: 14
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75009|AVENUE Count: 14
```

Figure 8: Groupement par Année d'autorisation, City|Codepostal|OpeningYear

```
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75014|BOULEVARD Count: 112
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75016|RUE Count: 518
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75017|RUE Count: 644
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75017|BOULEVARD Count: 154
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75017|AVENUE Count: 308
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75019|AVENUE Count: 252
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|PostalCode|StreetType: PARIS|75019|RUE Count: 336
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|StreetType|OpeningDate|AuthorizationYear: PARIS|PL|1994|1942 Count: 2
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|StreetType|OpeningDate|AuthorizationYear: PARIS|AVENUE|2008|2007 Count: 2
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|StreetType|OpeningDate|AuthorizationYear: PARIS|RUE|2007|1984 Count: 2
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|StreetType|OpeningDate|AuthorizationYear: PARIS|BOULEVARD|2011|2003 Count: 2
2025-01-20 13:30:13 INFO t.f.n.p.s.PharmaciesStreamProcessor - City|StreetType|OpeningDate|AuthorizationYear: PARIS|RUE|1985|1945 Count: 2
```

Figure 9: Groupement par City|PostalCode|StreetType et Groupement par City|StreetType|OpeningDate|AuthorizationYear

Partition	Offset	Key	Value	Timestamp
0	0	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:48.857
1	1	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:48.178
2	2	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:48.821
3	3	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:49.911
4	4	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:49.849
5	5	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.162
6	6	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.218
7	7	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.279
8	8	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.314
9	9	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.318
10	10	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.378
11	11	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.440
12	12	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.534
13	13	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.566
14	14	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.601
15	15	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.616
16	16	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.632
17	17	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.638
18	18	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.678
19	19	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.683
20	20	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.692
21	21	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.700
22	22	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.739
23	23	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.745
24	24	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.807
25	25	50415249537C3735303137C504C	00000000000000000000	2025-01-19 16:19:50.837

Figure 10: Affichage des Topics dans OffsetExplorer

Conclusion et Bibliographie

Ce projet illustre l'utilisation d'une architecture Kafka pour traiter, analyser et exposer des données. Grâce à une séparation des responsabilités entre les différentes applications (Producer, Consumer et Streams), chaque composant joue un rôle précis dans le pipeline. Le prétraitement des données, leur stockage dans une base MySQL et l'exposition via une API REST sont réalisés dans ce TP. Enfin, l'intégration d'outils comme Postman permet de tester les services développés.

<https://www.data.gouv.fr/fr/datasets/carte-des-pharmacies-de-paris/>

<https://github.com/loicosquare/TP-NFE-101.git>

<https://github.com/features/copilot>