

TP3 – Data Lake & Data Warehouse (Kafka, KSQLDB, Python, SQLite/MySQL)

1. Design du Data Lake

Élément	Description
Emplacement	Le Data Lake est un dossier local nommé <code>data_lake/</code> sur la machine.
Structure	Les données provenant de Kafka/KSQLDB sont stockées dans des sous-dossiers par topic et date : <code>data_lake/<topic>/<YYYY-MM-DD>/<topic>_<YYYY-MM-DD>.jsonl</code>
Partitionnement	Partition par date (ou version) pour faciliter la traçabilité et les chargements incrémentaux.
Ajout de nouveaux feeds	Lorsqu'un nouveau flux Kafka est ajouté : 1) Déclaration dans <code>config/settings.py</code> 2) Le consumer Kafka crée automatiquement le dossier ; 3) Le pipeline est réutilisé sans modification du code.
Stockage Streams & Tables	Chaque stream et table est stocké sous forme de fichiers JSONL distincts.
Mode de stockage	Streams → append ; Tables agrégées → overwrite ; Données statiques → ignore.

Structure

```
```text
```

```
data_lake/
```

```
|— config/
```

```
| |— settings.py
```

```
|— consumers/
```

```
| |— kafka_to_datalake.py
```

```
| |— fake_kafka_producer.py
```

```
|— jobs/
```

```
| |— sqlite_loader.py
```

```
|— pipeline/
```

```
| |— TRANSACTIONS_SECURE/
```

```
| |— TRANSACTIONS_USD/
```

```
| |— TRANSACTIONS_BLACKLISTED/
```

```
| |— USER_SPENDING_BY_TYPE/
```

```
| |— SPEND_LAST_FIVE_MIN_BY_TYPE/
```

```
|— data_warehouse.db
```

```
```
```

2. Design du Data Warehouse (MySQL / SQLite)

Toutes les tables de KSQLDB sont répliquées dans le Data Warehouse via un job Python (sqlite_loader.py).

| Table | Clé primaire | Clés étrangères | Description |
|------------------------------------|-----------------------------|-----------------|-------------------------------------|
| TRANSACTIONS_SECURE | transaction_id | user_id | Transactions brutes anonymisées |
| TRANSACTIONS_USD | transaction_id | user_id | Transactions converties en USD |
| USER_SPENDING_BY_TYPE | (user_id, transaction_type) | user_id | Agrégations par type et utilisateur |
| SPEND_LAST_FIVE_MIN_BY_TYPE | transaction_type | - | Dépense glissante 5 min |
| USERS | user_id | - | Utilisateurs anonymisés |
| PERMISSIONS | (user_id, topic_name) | user_id | Gestion des accès par utilisateur |

Diagramme Entité-Association (textuel) :

USERS (user_id PK)

- |
- |< TRANSACTIONS_SECURE (transaction_id PK, user_id FK)
- |< TRANSACTIONS_USD (transaction_id PK, user_id FK)
- |< USER_SPENDING_BY_TYPE (user_id FK, transaction_type, total_spent_usd)

PERMISSIONS (user_id FK, topic_name, access_level)

```
Windows PowerShell
(venv) PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\Datalakes_&_Data_Integration\TP\data_lake> python -m jobs.sqlite_loader

Tous les dossiers du Data Lake ont été vérifiés/cr  s (2025-10-31).
[  ] Donn  es import  es dans table transactions_secure depuis pipeline\transactions_secure\2025-10-31
[  ] Donn  es import  es dans table transactions_usd depuis pipeline\transactions_usd\2025-10-31
[  ] Donn  es import  es dans table transactions_blacklisted depuis pipeline\transactions_blacklisted\2025-10-31
[  ] Donn  es import  es dans table transactions_completed depuis pipeline\transactions_completed\2025-10-31
[  ] Donn  es import  es dans table transactions_failed depuis pipeline\transactions_failed\2025-10-31
[  ] Donn  es import  es dans table transactions_pending depuis pipeline\transactions_pending\2025-10-31
[  ] Donn  es import  es dans table transactions_processing depuis pipeline\transactions_processing\2025-10-31
[  ] Donn  es import  es dans table transactions_cancelled depuis pipeline\transactions_cancelled\2025-10-31
[  ] Donn  es import  es dans table spend_last_five_min_by_type depuis pipeline\spend_last_five_min_by_type\2025-10-31
[  ] Donn  es import  es dans table user_spending_by_type depuis pipeline\user_spending_by_type\2025-10-31
[  ] Import Data Lake -> SQLite termin  .
(venv) PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\Datalakes_&_Data_Integration\TP\data_lake>
```

data_lake

Version control

main

Start Free Trial

TRANSACTIONS_CANCELLED_2025-10-31.jsonl | sqlite_loader.py | data_warehouse.db | settings.py | requirements.txt | scheduler.py | sqlite_test_loader.py | logger.py

Tables Database Metadata

Table: transactions_usd Page: 0 Jump << < 1-22 > >> Refresh

| id | data | importe... |
|----|------------|-------------|
| 34 | { "TRAN... | 2025-10-... |
| 35 | { "TRAN... | 2025-10-... |
| 36 | { "TRAN... | 2025-10-... |
| 37 | { "TRAN... | 2025-10-... |
| 38 | { "TRAN... | 2025-10-... |
| 39 | { "TRAN... | 2025-10-... |
| 40 | { "TRAN... | 2025-10-... |
| 41 | { "TRAN... | 2025-10-... |
| 42 | { "TRAN... | 2025-10-... |
| 43 | { "TRAN... | 2025-10-... |
| 44 | { "TRAN... | 2025-10-... |
| 45 | { "TRAN... | 2025-10-... |
| 46 | { "TRAN... | 2025-10-... |
| 47 | { "TRAN... | 2025-10-... |
| 48 | { "TRAN... | 2025-10-... |
| 49 | { "TRAN... | 2025-10-... |
| 50 | { "TRAN... | 2025-10-... |

Transaction details for row 35:

{ "TRANSACTION_ID": "TXN-0bc9a602", "USER_ID_HASH": "7ffab80547e38b4308ddef923f4efb83", "USER_NAME": "E_HASHED": "660162d573b358f079be450220c015de", "TRANSACTION_TYPE": null, "AMOUNT": 923.91, "CURRENCY": "JPY", "CITY": "Chennai", "COUNTRY": "India", "IP_MASKED": "2.51.237.50", "STREET_MASKED": "54 Main St", "AMOUNT_USD": 6.190197, "TARGET_CURRENCY": "USD" }

data_lake > data_warehouse.db Python 3.11

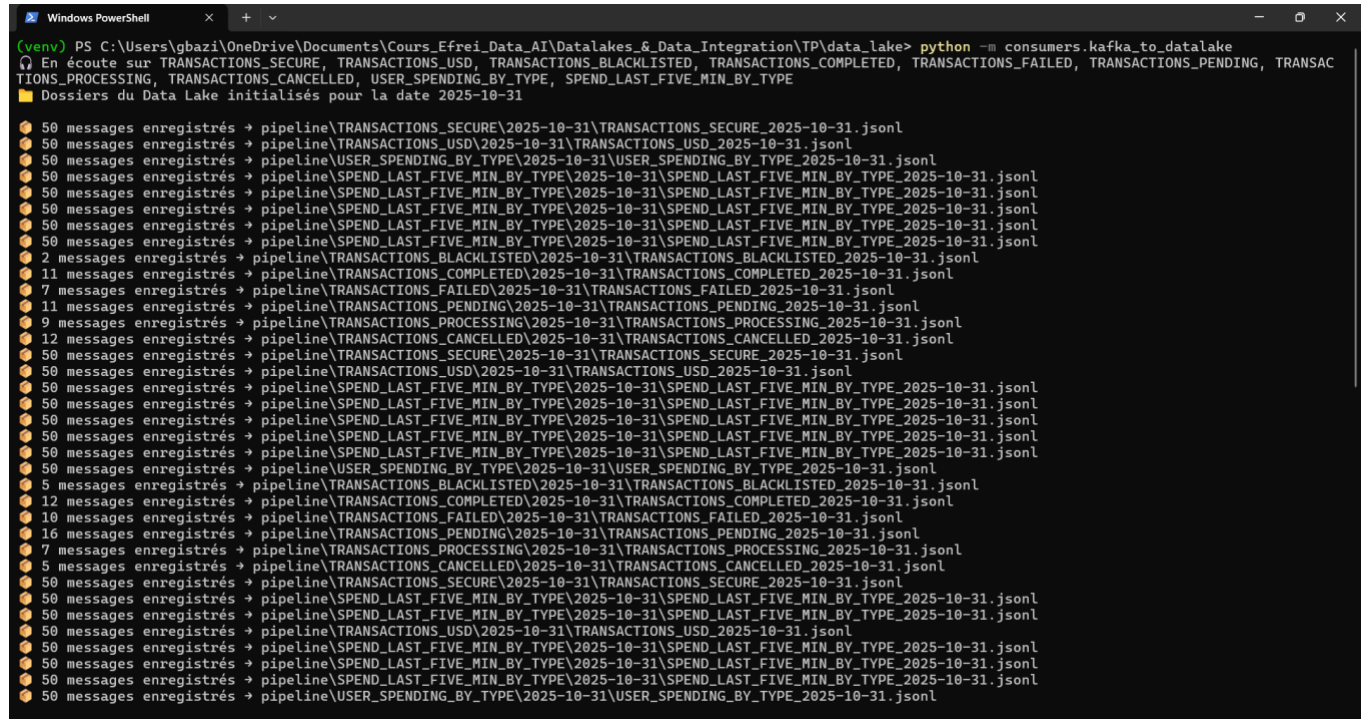
3. Kafka Consumers

Les Kafka Consumers permettent de lire les messages, les écrire dans le Data Lake et charger les données dans le Data Warehouse.

Exemple simplifié :

```
from kafka import KafkaConsumer
import json, os, datetime

consumer = KafkaConsumer('TRANSACTIONS_SECURE',
bootstrap_servers='localhost:9092')
for msg in consumer:
    data = json.loads(msg.value)
    date_str = datetime.date.today().isoformat()
    path = f'data_lake/TRANSACTIONS_SECURE/{date_str}/'
    os.makedirs(path, exist_ok=True)
    with open(f'{path}/data.jsonl', "a") as f:
        f.write(json.dumps(data) + "\n")
```



```
Windows PowerShell
PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\Datalakes & Data Integration\TP\data_lake> python -m consumers.kafka_to_dataLake
En écoute sur TRANSACTIONS_SECURE, TRANSACTIONS_USD, TRANSACTIONS_BLACKLISTED, TRANSACTIONS_COMPLETED, TRANSACTIONS_FAILED, TRANSACTIONS_PENDING, TRANSACTIONS_PROCESSING, TRANSACTIONS_CANCELLED, USER_SPENDING_BY_TYPE, SPEND_LAST_FIVE_MIN_BY_TYPE
Dossiers du Data Lake initialisés pour la date 2025-10-31
50 messages enregistrés -> pipeline\TRANSACTIONS_SECURE\2025-10-31\TRANSACTIONS_SECURE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\TRANSACTIONS_USD\2025-10-31\TRANSACTIONS_USD_2025-10-31.jsonl
50 messages enregistrés -> pipeline\USER_SPENDING_BY_TYPE\2025-10-31\USER_SPENDING_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\TRANSACTIONS_BLACKLISTED\2025-10-31\TRANSACTIONS_BLACKLISTED_2025-10-31.jsonl
2 messages enregistrés -> pipeline\TRANSACTIONS_BLACKLISTED\2025-10-31\TRANSACTIONS_BLACKLISTED_2025-10-31.jsonl
11 messages enregistrés -> pipeline\TRANSACTIONS_COMPLETED\2025-10-31\TRANSACTIONS_COMPLETED_2025-10-31.jsonl
7 messages enregistrés -> pipeline\TRANSACTIONS_FAILED\2025-10-31\TRANSACTIONS_FAILED_2025-10-31.jsonl
11 messages enregistrés -> pipeline\TRANSACTIONS_PENDING\2025-10-31\TRANSACTIONS_PENDING_2025-10-31.jsonl
9 messages enregistrés -> pipeline\TRANSACTIONS_PROCESSING\2025-10-31\TRANSACTIONS_PROCESSING_2025-10-31.jsonl
12 messages enregistrés -> pipeline\TRANSACTIONS_CANCELLED\2025-10-31\TRANSACTIONS_CANCELLED_2025-10-31.jsonl
50 messages enregistrés -> pipeline\TRANSACTIONS_SECURE\2025-10-31\TRANSACTIONS_SECURE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\TRANSACTIONS_USD\2025-10-31\TRANSACTIONS_USD_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\USER_SPENDING_BY_TYPE\2025-10-31\USER_SPENDING_BY_TYPE_2025-10-31.jsonl
5 messages enregistrés -> pipeline\TRANSACTIONS_BLACKLISTED\2025-10-31\TRANSACTIONS_BLACKLISTED_2025-10-31.jsonl
12 messages enregistrés -> pipeline\TRANSACTIONS_COMPLETED\2025-10-31\TRANSACTIONS_COMPLETED_2025-10-31.jsonl
10 messages enregistrés -> pipeline\TRANSACTIONS_FAILED\2025-10-31\TRANSACTIONS_FAILED_2025-10-31.jsonl
16 messages enregistrés -> pipeline\TRANSACTIONS_PENDING\2025-10-31\TRANSACTIONS_PENDING_2025-10-31.jsonl
7 messages enregistrés -> pipeline\TRANSACTIONS_PROCESSING\2025-10-31\TRANSACTIONS_PROCESSING_2025-10-31.jsonl
5 messages enregistrés -> pipeline\TRANSACTIONS_CANCELLED\2025-10-31\TRANSACTIONS_CANCELLED_2025-10-31.jsonl
50 messages enregistrés -> pipeline\TRANSACTIONS_SECURE\2025-10-31\TRANSACTIONS_SECURE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\TRANSACTIONS_USD\2025-10-31\TRANSACTIONS_USD_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés -> pipeline\USER_SPENDING_BY_TYPE\2025-10-31\USER_SPENDING_BY_TYPE_2025-10-31.jsonl
```

4. Gouvernance et Sécurité

| Domaine | Mécanisme | Objectif |
|------------------------|--|----------------------------------|
| Suppression historique | Script Python supprimant les données après 30 jours | Optimisation & conformité RGPD |
| Gestion des droits | Table permissions (user_id, topic_name, access_level) | Contrôle des accès |
| Feed Kafka | Déclaration dans settings.py → dossier auto → pipeline réutilisé | Gain de temps |
| Anonymisation | Masquage KSQLDB (MASK, HASH, CONCAT) | Protection des données sensibles |
| Traçabilité | Transaction_id + timestamp | Audit & transparence |
| | | |

```
Windows PowerShell
(venv) PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\Datalakes_&_Data_Integration\TP\data_lake> python -m job
s.cleanup
(venv) PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\Datalakes_&_Data_Integration\TP\data_lake> python -m job
s.permissions_manager
✅ Table 'permissions' initialisée.
🔔 Permission ajoutée pour admin sur TRANSACTIONS_SECURE.
🔔 Permission ajoutée pour analyst sur TRANSACTIONS_USD.
[[('1', 'admin', 'TRANSACTIONS_SECURE', 1, 1, '2025-11-01 19:24:28'), ('2', 'analyst', 'TRANSACTIONS_USD', 1, 0, '2025-11-01
19:24:28')]]
(venv) PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\Datalakes_&_Data_Integration\TP\data_lake>
```

5. Orchestration et Optimisation

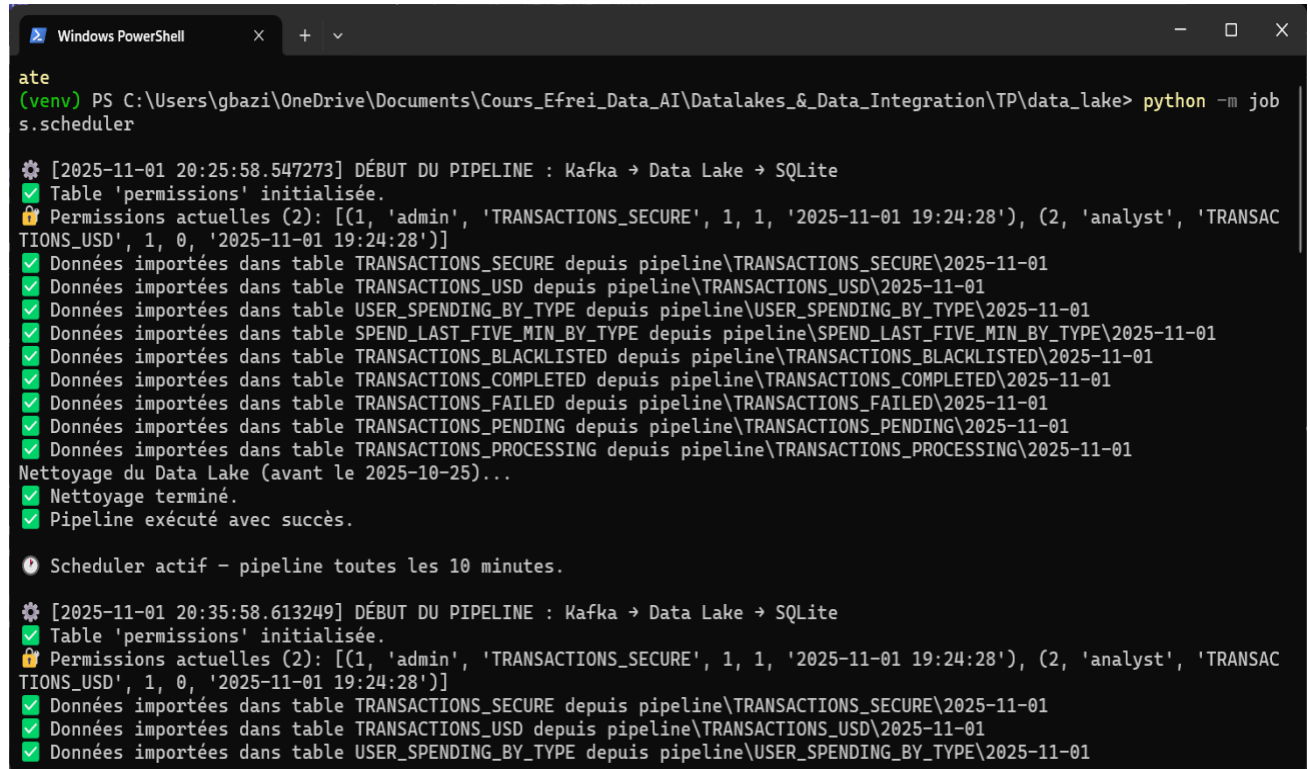
Orchestration automatisée via schedule ou Apache Beam (toutes les 10 minutes).

```
import schedule, os, time
```

```
def run_jobs():  
    os.system("python consumers/kafka_to_datalake.py")  
    os.system("python jobs/sqlite_loader.py")
```

```
schedule.every(10).minutes.do(run_jobs)
```

```
while True:  
    schedule.run_pending()  
    time.sleep(60)
```



```
ate  
(venv) PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\Datalakes_&_Data_Integration\TP\data_lake> python -m jobs.scheduler  
  
[2025-11-01 20:25:58.547273] DÉBUT DU PIPELINE : Kafka → Data Lake → SQLite  
✓ Table 'permissions' initialisée.  
🔑 Permissions actuelles (2): [(1, 'admin', 'TRANSACTIONS_SECURE', 1, 1, '2025-11-01 19:24:28'), (2, 'analyst', 'TRANSACTIONS_USD', 1, 0, '2025-11-01 19:24:28')]  
✓ Données importées dans table TRANSACTIONS_SECURE depuis pipeline\TRANSACTIONS_SECURE\2025-11-01  
✓ Données importées dans table TRANSACTIONS_USD depuis pipeline\TRANSACTIONS_USD\2025-11-01  
✓ Données importées dans table USER_SPENDING_BY_TYPE depuis pipeline\USER_SPENDING_BY_TYPE\2025-11-01  
✓ Données importées dans table SPEND_LAST_FIVE_MIN_BY_TYPE depuis pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-11-01  
✓ Données importées dans table TRANSACTIONS_BLACKLISTED depuis pipeline\TRANSACTIONS_BLACKLISTED\2025-11-01  
✓ Données importées dans table TRANSACTIONS_COMPLETED depuis pipeline\TRANSACTIONS_COMPLETED\2025-11-01  
✓ Données importées dans table TRANSACTIONS_FAILED depuis pipeline\TRANSACTIONS_FAILED\2025-11-01  
✓ Données importées dans table TRANSACTIONS_PENDING depuis pipeline\TRANSACTIONS_PENDING\2025-11-01  
✓ Données importées dans table TRANSACTIONS_PROCESSING depuis pipeline\TRANSACTIONS_PROCESSING\2025-11-01  
Nettoyage du Data Lake (avant le 2025-10-25)...  
✓ Nettoyage terminé.  
✓ Pipeline exécuté avec succès.  
  
🕒 Scheduler actif - pipeline toutes les 10 minutes.  
  
[2025-11-01 20:35:58.613249] DÉBUT DU PIPELINE : Kafka → Data Lake → SQLite  
✓ Table 'permissions' initialisée.  
🔑 Permissions actuelles (2): [(1, 'admin', 'TRANSACTIONS_SECURE', 1, 1, '2025-11-01 19:24:28'), (2, 'analyst', 'TRANSACTIONS_USD', 1, 0, '2025-11-01 19:24:28')]  
✓ Données importées dans table TRANSACTIONS_SECURE depuis pipeline\TRANSACTIONS_SECURE\2025-11-01  
✓ Données importées dans table TRANSACTIONS_USD depuis pipeline\TRANSACTIONS_USD\2025-11-01  
✓ Données importées dans table USER_SPENDING_BY_TYPE depuis pipeline\USER_SPENDING_BY_TYPE\2025-11-01
```

Synthèse

| Domaine | Résumé |
|------------------------|---|
| Data Lake | Stockage brut JSON partitionné par date, extensible, append/overwrite selon le type de données. |
| Data Warehouse | Schéma relationnel clair (SQLite/MySQL), clés primaires et étrangères définies. |
| Kafka Consumers | Automatisation de la collecte et intégration continue. |
| Sécurité & Gouvernance | Rétention, anonymisation, permissions et logs d'audit. |
| Orchestration | Automatisation des jobs toutes les 10 minutes via schedule / Apache Beam. |

```
Windows PowerShell
PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\Datalakes_&_Data_Integration\TP\data_lake> .\venv\Scripts\activate
(venv) PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\Datalakes_&_Data_Integration\TP\data_lake> python -m jobs.scheduler

[2025-11-01 20:25:58.547273] DÉBUT DU PIPELINE : Kafka → Data Lake → SQLite
✓ Table 'permissions' initialisée.
🔔 Permissions actuelles (2): [(1, 'admin', 'TRANSACTIONS_SECURE', 1, 1, '2025-11-01 19:24:28'), (2, 'analyst', 'TRANSACTIONS_USD', 1, 0, '2025-11-01 19:24:28')]
✓ Données importées dans table TRANSACTIONS_SECURE depuis pipeline\TRANSACTIONS_SECURE\2025-11-01
✓ Données importées dans table TRANSACTIONS_USD depuis pipeline\TRANSACTIONS_USD\2025-11-01
✓ Données importées dans table USER_SPENDING_BY_TYPE depuis pipeline\USER_SPENDING_BY_TYPE\2025-11-01
✓ Données importées dans table SPEND_LAST_FIVE_MIN_BY_TYPE depuis pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-11-01
✓ Données importées dans table TRANSACTIONS_BLACKLISTED depuis pipeline\TRANSACTIONS_BLACKLISTED\2025-11-01
✓ Données importées dans table TRANSACTIONS_COMPLETED depuis pipeline\TRANSACTIONS_COMPLETED\2025-11-01
✓ Données importées dans table TRANSACTIONS_FAILED depuis pipeline\TRANSACTIONS_FAILED\2025-11-01
✓ Données importées dans table TRANSACTIONS_PENDING depuis pipeline\TRANSACTIONS_PENDING\2025-11-01
✓ Données importées dans table TRANSACTIONS_PROCESSING depuis pipeline\TRANSACTIONS_PROCESSING\2025-11-01
Nettoyage du Data Lake (avant le 2025-10-25)...
✓ Nettoyage terminé.
✓ Pipeline exécuté avec succès.

🕒 Scheduler actif - pipeline toutes les 10 minutes.
```



```
Windows PowerShell
(venv) PS C:\Users\gbazi\OneDrive\Documents\Cours_Efrei_Data_AI\DataLakes_&_Data_Integration\TP\data_lake> python -m consumers.kafka_to_datalake
En écoute sur TRANSACTIONS_SECURE, TRANSACTIONS_USD, TRANSACTIONS_BLACKLISTED, TRANSACTIONS_COMPLETED, TRANSACTIONS_FAILED, TRANSACTIONS_PENDING, TRANSACTIONS_PROCESSING, TRANSACTIONS_CANCELLED, USER_SPENDING_BY_TYPE, SPEND_LAST_FIVE_MIN_BY_TYPE
Dossiers du Data Lake initialisés pour la date 2025-10-31

50 messages enregistrés + pipeline\TRANSACTIONS_SECURE\2025-10-31\TRANSACTIONS_SECURE_2025-10-31.jsonl
50 messages enregistrés + pipeline\TRANSACTIONS_USD\2025-10-31\TRANSACTIONS_USD_2025-10-31.jsonl
50 messages enregistrés + pipeline\USER_SPENDING_BY_TYPE\2025-10-31\USER_SPENDING_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
2 messages enregistrés + pipeline\TRANSACTIONS_BLACKLISTED\2025-10-31\TRANSACTIONS_BLACKLISTED_2025-10-31.jsonl
11 messages enregistrés + pipeline\TRANSACTIONS_COMPLETED\2025-10-31\TRANSACTIONS_COMPLETED_2025-10-31.jsonl
7 messages enregistrés + pipeline\TRANSACTIONS_FAILED\2025-10-31\TRANSACTIONS_FAILED_2025-10-31.jsonl
11 messages enregistrés + pipeline\TRANSACTIONS_PENDING\2025-10-31\TRANSACTIONS_PENDING_2025-10-31.jsonl
9 messages enregistrés + pipeline\TRANSACTIONS_PROCESSING\2025-10-31\TRANSACTIONS_PROCESSING_2025-10-31.jsonl
12 messages enregistrés + pipeline\TRANSACTIONS_CANCELLED\2025-10-31\TRANSACTIONS_CANCELLED_2025-10-31.jsonl
50 messages enregistrés + pipeline\TRANSACTIONS_SECURE\2025-10-31\TRANSACTIONS_SECURE_2025-10-31.jsonl
50 messages enregistrés + pipeline\TRANSACTIONS_USD\2025-10-31\TRANSACTIONS_USD_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
5 messages enregistrés + pipeline\TRANSACTIONS_BLACKLISTED\2025-10-31\TRANSACTIONS_BLACKLISTED_2025-10-31.jsonl
12 messages enregistrés + pipeline\TRANSACTIONS_COMPLETED\2025-10-31\TRANSACTIONS_COMPLETED_2025-10-31.jsonl
10 messages enregistrés + pipeline\TRANSACTIONS_FAILED\2025-10-31\TRANSACTIONS_FAILED_2025-10-31.jsonl
16 messages enregistrés + pipeline\TRANSACTIONS_PENDING\2025-10-31\TRANSACTIONS_PENDING_2025-10-31.jsonl
7 messages enregistrés + pipeline\TRANSACTIONS_PROCESSING\2025-10-31\TRANSACTIONS_PROCESSING_2025-10-31.jsonl
5 messages enregistrés + pipeline\TRANSACTIONS_CANCELLED\2025-10-31\TRANSACTIONS_CANCELLED_2025-10-31.jsonl
50 messages enregistrés + pipeline\TRANSACTIONS_SECURE\2025-10-31\TRANSACTIONS_SECURE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\TRANSACTIONS_USD\2025-10-31\TRANSACTIONS_USD_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\SPEND_LAST_FIVE_MIN_BY_TYPE\2025-10-31\SPEND_LAST_FIVE_MIN_BY_TYPE_2025-10-31.jsonl
50 messages enregistrés + pipeline\USER_SPENDING_BY_TYPE\2025-10-31\USER_SPENDING_BY_TYPE_2025-10-31.jsonl
```