

**DATA SCIENCE & ANALYSE DE
DONNEES**



**DIFFUSION EN CONTINU
STRUCTURÉE AVEC APACHE SPARK**



Proposé par : _____

M. MOUNGOLE ARISTIDE

Réalisé par : _____

HYONTA KENGAP BLEROT

Année : 2024/2025

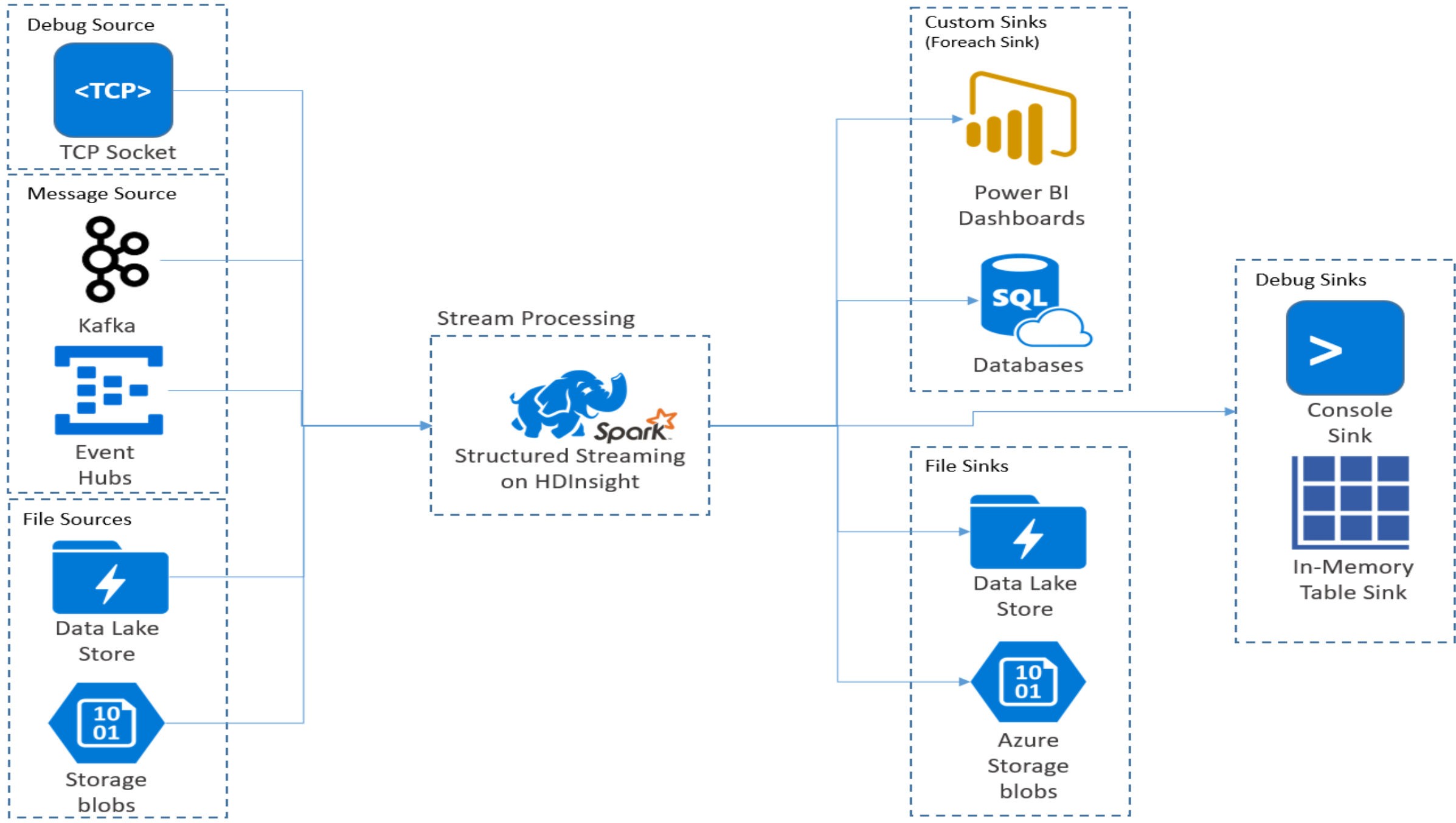


Table des matières

Introduction à la Diffusion en Continu Structurée

01

Qu'est-ce que la Diffusion en Continu Structurée ?

02

Technologies Utilisées

03

Contexte Camerounais et Applications

04

Secteur Financier : Détection de Fraudes

05

Architecture Globale du Système

06

Exemple de Code : Producteur Kafka en Python

07

Traitement avec Apache Spark en Scala

08

Démonstration : Détection de Fraude en Temps Réel

09

Conclusion

10

Introduction à la Diffusion en Continu Structurée



Introduction à la Diffusion en Continu Structurée

La diffusion en continu structurée est une approche moderne du traitement des flux de données en temps réel.

Cette approche gère efficacement les données de différentes sources.



Applications Potentielles

Dans le domaine de l'éducation et de la formation, elle peut être utilisée pour fournir des contenus d'apprentissage en temps réel.

En marketing et communication, elle permet d'analyser les comportements des consommateurs et d'ajuster les stratégies en conséquence.

Dans la recherche et développement, elle facilite l'accès aux données et aux résultats d'expériences en temps réel.



Importance de la Diffusion

Elle facilite la circulation des idées et des informations entre les utilisateurs.

Elle renforce la collaboration entre les individus et les organisations, favorisant ainsi un environnement de travail dynamique.

Elle contribue à l'innovation et à l'amélioration des processus en permettant une analyse rapide des données.



Conclusion

La diffusion en continu structurée représente une avancée significative dans le traitement des données, offrant des solutions adaptées à divers secteurs.

Son adoption peut transformer la manière dont les organisations gèrent et analysent les informations.

Qu'est-ce que la Diffusion en Continu Structurée ?



Définition de la diffusion en continu structurée

La diffusion en continu structurée est un moteur de traitement de flux de données en temps réel introduit dans Apache Spark à partir de la version 2.0.

Elle repose sur un modèle de programmation déclaratif utilisant Spark SQL.



Comparaison avec d'autres technologies de traitement de flux

Contrairement aux approches traditionnelles comme Apache Storm ou Flink, la diffusion en continu structurée utilise un modèle de micro-lots.

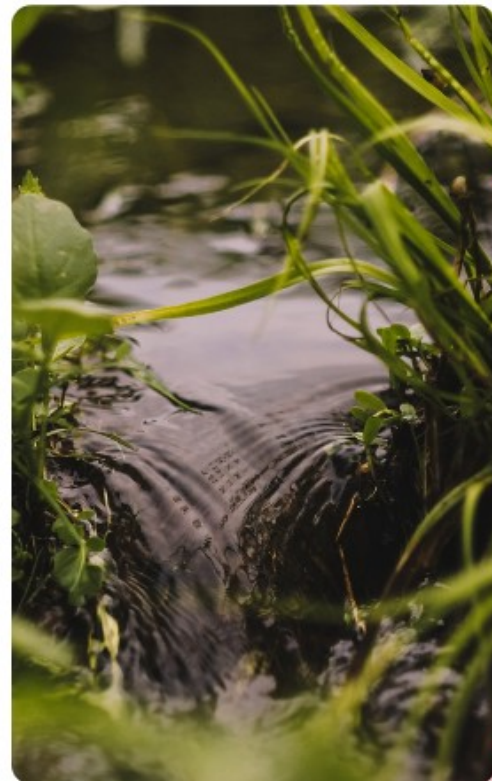
Cette méthode permet une gestion efficace des états et une intégration native avec l'écosystème Spark.



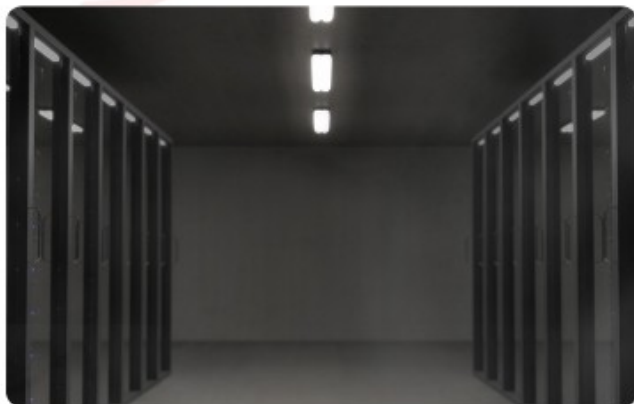
Avantages spécifiques de l'approche

Elle permet de traiter les flux de données de manière similaire aux bases de données relationnelles.

La diffusion en continu structurée est particulièrement adaptée pour des applications nécessitant une détection de fraudes en temps réel.



Technologies Utilisées



Technologies Utilisées

Présentation des technologies utilisées dans le projet.



Apache Spark

Apache Spark est un moteur de traitement de flux utilisé pour le traitement en temps réel des données.



Apache Kafka

Apache Kafka est un système de messagerie distribué pour la diffusion des données en temps réel.



PostgreSQL

PostgreSQL est une base de données relationnelle pour stocker les résultats du traitement.



Scala

Scala est le langage de programmation utilisé pour développer l'application Spark.

Contexte Camerounais et Applications



Contexte Camerounais et Applications

La diffusion en continu structurée est essentielle pour le traitement des données en temps réel au Cameroun.

Elle a des applications significatives dans le secteur bancaire, notamment pour la détection de fraudes.

Secteur Financier : Détection de Fraudes



Mécanisme de Détection de Fraudes

Les transactions dépassant un seuil de 10 000 XAF sont considérées comme suspectes et stockées pour une analyse approfondie.

Ce système permet aux banques camerounaises de prendre des mesures immédiates pour enquêter et éventuellement bloquer une transaction.



Traitement en Temps Réel

L'utilisation d'Apache Spark permet un traitement en temps réel des données, essentiel pour améliorer la sécurité et l'efficacité des systèmes financiers.

Les données sont analysées instantanément pour détecter des activités frauduleuses.



Surveillance des Transactions

La diffusion en continu structurée permet de surveiller les transactions en temps réel et d'identifier des anomalies.

Des critères spécifiques, tels que la fréquence anormale et la localisation suspecte, sont utilisés pour détecter les fraudes.

Architecture Globale du Système



Exemple de Code : Producteur Kafka en Python



Exemple de Code : Producteur Kafka en Python

Ce script envoie des transactions à Kafka.

```
from kafka import KafkaProducer
import json
import time
import random

def generate_transaction():
    return {
        'id': random.randint(1000, 9999),
        'transaction_id': random.randint(10000, 99999),
        'amount': round(random.uniform(100, 20000), 2),
        'currency': 'XAF',
        'timestamp': time.strftime('%Y-%m-%d %H:%M:%S')
    }

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

while True:
    transaction = generate_transaction()
    producer.send('transaction', transaction)
    print(f'Transaction envoyée : {transaction}')
    time.sleep(1)
```

Traitement avec Apache Spark en Scala



Filtrage des Transactions en Temps Réel

Le script lit les flux de transactions à partir de Kafka.

Les transactions dont le montant dépasse 10 000 XAF sont filtrées.



Logique de Détection de Fraude

Les transactions suspectes sont stockées pour une analyse approfondie.

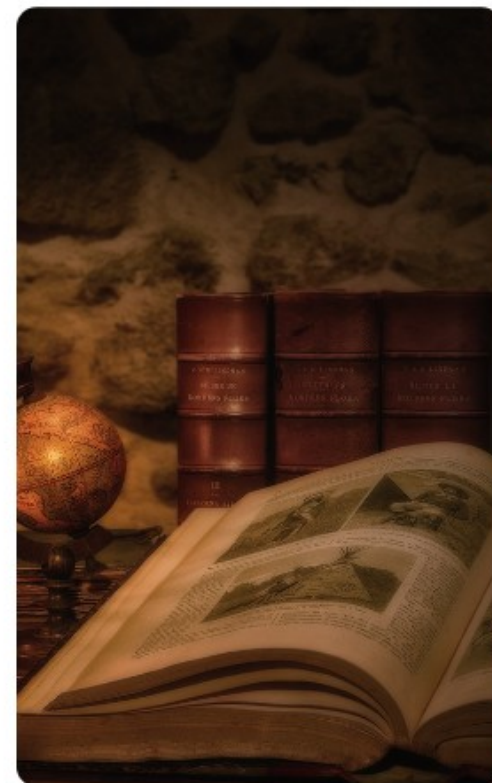
Le système permet aux banques de prendre des mesures immédiates.



Stockage des Données dans PostgreSQL

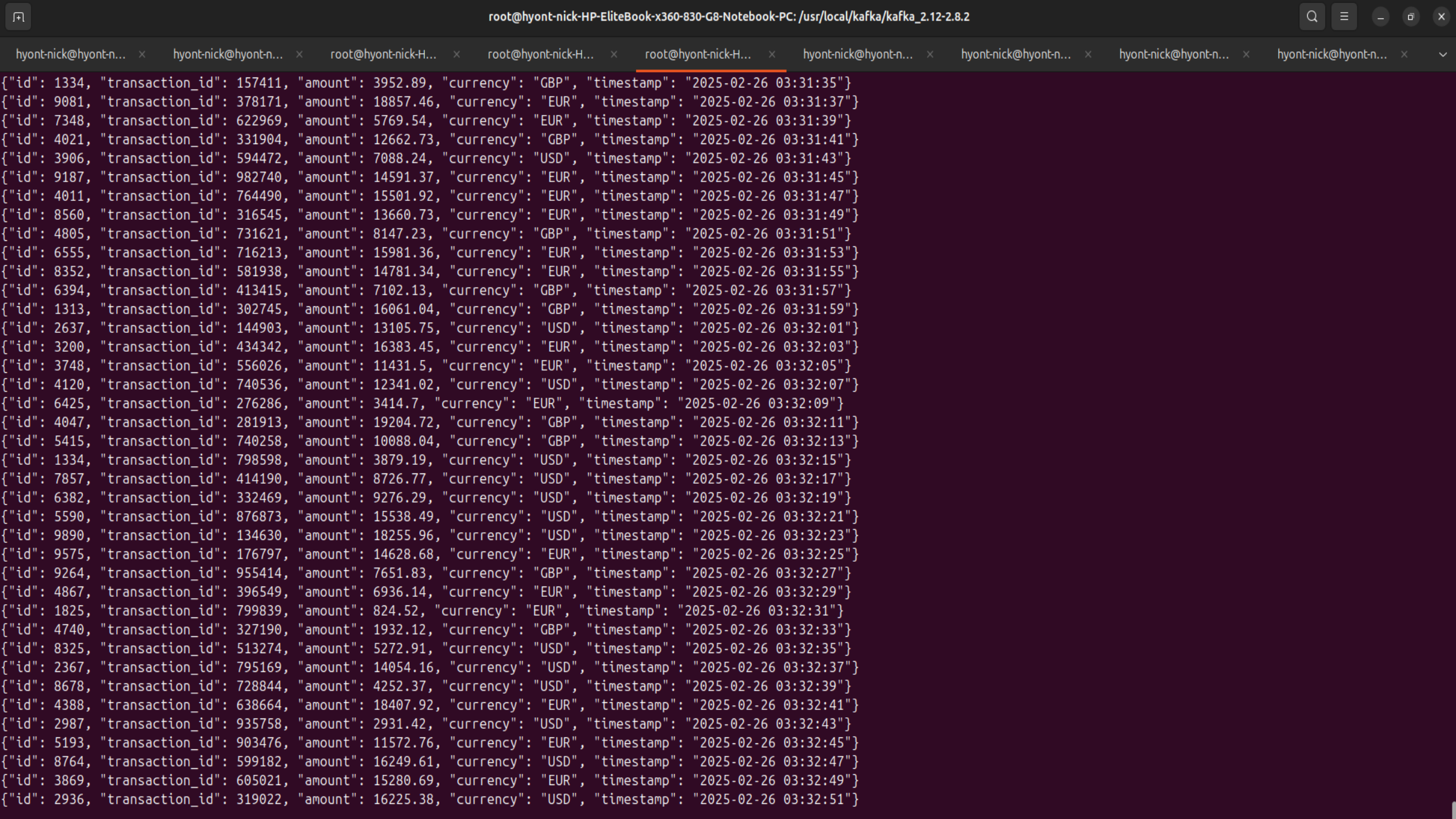
Les transactions filtrées sont écrites dans une base de données PostgreSQL.

Cela facilite le suivi et l'analyse des transactions suspectes.



Démonstration : Détection de Fraude en Temps Réel





Conclusion

Résumé des points clés du projet

La diffusion en continu structurée avec Apache Spark permet de traiter des flux de données en temps réel, en particulier pour la détection de fraudes dans le secteur bancaire camerounais.

Avantages de la diffusion en continu structurée

Cette technologie offre une gestion efficace des données, permettant aux banques de détecter immédiatement des activités suspectes et de renforcer la sécurité des transactions.

Perspectives d'avenir

Avec l'augmentation des transactions numériques, l'adoption de solutions de traitement en temps réel est essentielle pour améliorer la sécurité et l'efficacité des systèmes financiers au Cameroun.

Merci pour votre attention

