**ABSTRACT: Fraud Detection in Banking Transactions Using Hadoop**

**Introduction**

This project presents a **Big Data solution** designed to process and analyse massive volumes of financial transaction data using the Hadoop ecosystem. By leveraging scalable distributed computing, this system helps financial institutions identify suspicious activities, detect fraudulent behaviour patterns, and enhance security and trust in digital banking. Real-time or batch-based fraud detection is made possible through intelligent rule-based checks and machine learning techniques built on top of Hadoop.

---

**Problem Statement**

Financial institutions face the challenge of processing enormous and continuous streams of transaction data from multiple sources. Traditional systems often struggle to analyse this data at scale and in real-time, leading to delays in fraud detection and potential financial losses. This project addresses the need for a **scalable and efficient system** that can ingest, store, and analyse vast datasets to identify fraudulent patterns and anomalies promptly, thereby bolstering security.

---

**Tools Used**

- **Data Storage**: Apache Hadoop HDFS

- **Data Processing**: MapReduce / Apache Spark

- **Data Querying**: Apache Hive / Pig

- **Data Streaming/Ingestion**: Apache Kafka / Flume

- **Data Import/Export**: Sqoop

- **Resource Management**: YARN

- **Machine Learning**: Apache Mahout, Spark MLlib, or other Python-based libraries

- **Visualization & Reporting**: Elasticsearch + Kibana / Tableau

- **Implementation Languages**: Java / Scala / Python

---

**System Modules**

1. **Data Ingestion Module**: Ingests large-scale banking transaction data from various sources, including real-time feeds from APIs, message queues, and logs using **Kafka** or **Flume**. It also handles bulk data import from traditional databases using **Sqoop**.

2. **Distributed Storage Module**: Securely stores structured and unstructured financial data, such as transaction records and customer profiles, in **HDFS**.

3. **Processing & Analytics Engine**: Uses **MapReduce** or **Spark** to analyze transaction history. It applies rule-based logic and machine learning models to detect anomalies, profile user behaviour, and validate transactions against blacklists and geolocation data.

4. **Alerting & Reporting Module**: Generates real-time alerts for flagged transactions and rule violations (e.g., unusually high transfers). It creates user behavior reports and aggregated insights on fraud trends, which are displayed on **Kibana** or **Tableau** dashboards for risk analysts.

---

**Flow**

The system's workflow follows a logical data pipeline:

**Data Sources** (RDBMS, APIs, Logs) → **Ingestion** (Sqoop, Kafka, Flume) → **Storage** (HDFS) → **Processing & Analysis** (Spark/MapReduce with ML Models) → **Outputs** (Alerts, Risk Scores, Reports) → **Visualization** (Dashboards on Kibana/Tableau)

---

**Conclusion**

This Hadoop-based fraud detection system offers significant benefits by enabling the **early detection of fraudulent transactions**, which prevents financial loss and supports compliance with AML and KYC regulations. It enhances trust in digital banking by automating anomaly detection at scale, reducing manual effort, and allowing the system to adapt to new fraud patterns through machine learning.

---

**Future Scope**

- **Advanced AI Integration**: Implement more sophisticated deep learning models (like LSTMs or Transformer networks) for more accurate sequential transaction analysis and anomaly detection.

- **Cloud Deployment**: Migrate the entire ecosystem to a cloud platform (like AWS, Azure, or GCP) to leverage managed services, improve scalability, and reduce infrastructure overhead.

- **Enhanced Real-Time Capabilities**: Integrate ultra-low-latency processing frameworks to decrease detection time from minutes to seconds.

<div align="center">

**24M11MC266---P.KUMARA SWAMY---ADITYA UNIVERSITY**

</div>