# PayPal fraud detection

One recent High Performance Data Analytics use case from the industry that gathered a lot of interest in the around the word is the the implementation of an HPDA system for real-time fraud detection across millions of daily transactions, successfully done by PayPal.
While internet commerce has become a vital part of the economy, detecting fraud in 'real time' as millions of
transactions are captured and processed by an assortment of systems – many having proprietary software
tools – has created a need to develop and deploy new fraud detection models.
PayPal, an eBay company, has used Hadoop and other software tools to detect fraud, but the colossal volumes of data were so large their systems were unable to perform the analysis quickly enough. To meet the challenge of finding fraud in near real time, PayPal decided to use HPC class systems – including the Lustre file system on their Hadoop cluster.
The result? In their first year of production PayPal saved over $700 million in fraudulent transactions that thy
would not have detected previously. [https://www.intel.com/content/dam/www/public/us/en/documents/white-papers/big-data-meets-high-performance-computing-white-paper.pdf]

## Core Infrastructure & Scale

PayPal's fraud detection system manages 74 petabytes of total storage with a 32% annual storage growth rate. [https://www.slideshare.net/slideshow/big-data-fast-data-paypal-yow-2018/126923757]The infrastructure processes 20 terabytes of log data collected daily and operates over 2,000 database instances handling 116 billion database calls per day. The system processes 41 million transactions per day, scaling up from 7.8 billion payment transactions annually in 2018. PayPal serves 400 million consumer accounts and 20 million merchant accounts across 200+ markets and currencies, processing a global transaction volume of $1.53 trillion in 2023.[https://helplama.com/paypal-revenue-users-statistics/,https://www.oberlo.com/statistics/us-paypal-transaction-volume]. Current Prevention stands for $500 million blocked in fraud per quarter.[https://chiefaiofficer.com/blog/how-paypals-ai-blocks-500-million-in-fraud-per-quarter/]

## Technical Architecture of PayPal's Fraud Detection System

PayPal's real-time fraud detection infrastructure is built on a sophisticated combination of big data and high-performance computing technologies designed to process massive transaction volumes with minimal latency. The system employs a Hadoop cluster configured with HBase for distributed data storage, enabling efficient management of petabyte-scale datasets across multiple nodes. This foundation provides the scalability and fault tolerance required for continuous operation across PayPal's global payment network.[https://intellipaat.com/blog/paypal-leverages-big-data-analytics/]

The processing layer utilizes a Lustre parallel file system running on top of Hadoop, combined with HPC-class infrastructure and high-speed InfiniBand interconnects. This architecture enables extremely fast data access and communication between compute nodes, which is critical for maintaining the millisecond-level latency

requirements of real-time fraud detection. The InfiniBand interconnects provide low-latency, high-bandwidth communication that significantly outperforms standard Ethernet networks, allowing the system to process fraud detection queries across thousands of transactions simultaneously without performance degradation.

For analytics and real-time streaming, PayPal leverages the Kraken analytics platform alongside Kafka for streaming data pipelines. Kafka enables the continuous ingestion of transaction data from multiple sources, while Kraken provides the analytics capabilities to derive insights and detect fraudulent patterns in real-time. This combination allows PayPal to process event streams of millions of transactions per day and apply complex fraud detection rules instantaneously. [https://www.vamsitalkstech.com/architecture/hadoop-counters-credit-card-fraud-23/]

The core database infrastructure consists of 2,000 Aerospike servers, including 200 high-density servers equipped with Intel Optane Persistent Memory. These Aerospike servers manage 100 petabytes of historical transaction and customer behavior data that inform fraud detection decisions.[https://pages.aerospike.com/rs/229-XUE-318/images/Aerospike_Case_Study_PayPal_Intel.pdf].

## Machine Learning Algorithms

- Algorithm Types: Neural Networks, Deep Learning, and Linear Regression

- Data Analysis: Thousands of data points including IP address, buying history, merchant activity, cookie information, device characteristics, geographic location, behavioral indicators

## Workflow Characteristics

- Continuous Operation: 24/7 real-time processing

- Human Review: Suspicious transactions flagged for hybrid automated and human expert scrutiny

- Adaptive Learning: Models update continuously based on new fraud patterns

## How the 5V Big Data Challenges Apply?

Volume: Manages 100 petabytes across infrastructure with 32% annual growth rate. Processes 41 million daily transactions generating 20+ terabytes of log data daily.

Velocity: Real-time processing with millisecond-level fraud detection on 41 million daily transactions. System handles 116 billion database calls per day, requiring near-instantaneous decisions to prevent fraud before transaction completion.

Variety: Heterogeneous data from "an assortment of systems – many having proprietary software tools", including structured transaction records, semi-structured logs, unstructured behavioral data, device fingerprints, network logs

spanning 200+ markets and currencies.

Veracity: 50% reduction in false positives while maintaining fraud detection effectiveness. Hybrid system combines machine learning (Neural Networks, Deep Learning, Linear Regression) with human detective verification to ensure accuracy.

Value: Directly quantifiable at $700 million fraud prevented in first year, now $500 million prevented quarterly. Enables PayPal to maintain trust across $1.53 trillion annual transaction volume

## Primary Citations:

1    SlideShare YOW 2018: "Big Data, Fast Data @ PayPal" - Infrastructure metrics

2    Intel White Paper 2014: "Big Data Meets High Performance Computing" - Original case study

3    Intel/Aerospike Case Study 2020: "PayPal Solves Fraud Challenges with Aerospike and Intel Optane PMem"

4    Various industry sources 2024-2025: Current performance metrics