

Operational Analytics of HDFS Logs with Anomaly Detection

Abhitosh - abhitosh2024@iisc.ac.in

Amit Nitin Joshi - amitj@iisc.ac.in

Naik Raghavendra Narottam - nraghavendra@iisc.ac.in

Vignesh S - vigneshs@iisc.ac.in

Problem statement

- Derive long term and real-time Operational Insights using Application Logs
 - System Health
 - System Performance
 - Anomalies
- Challenges
 - Scale and Volume of
 - Velocity in scaled environment

Features

- Operations dashboard
 - Real time insights
 - Daily insights
 - Blocks with anomalies
- Non Functional
 - Optimal throughput and latency

Design Goals

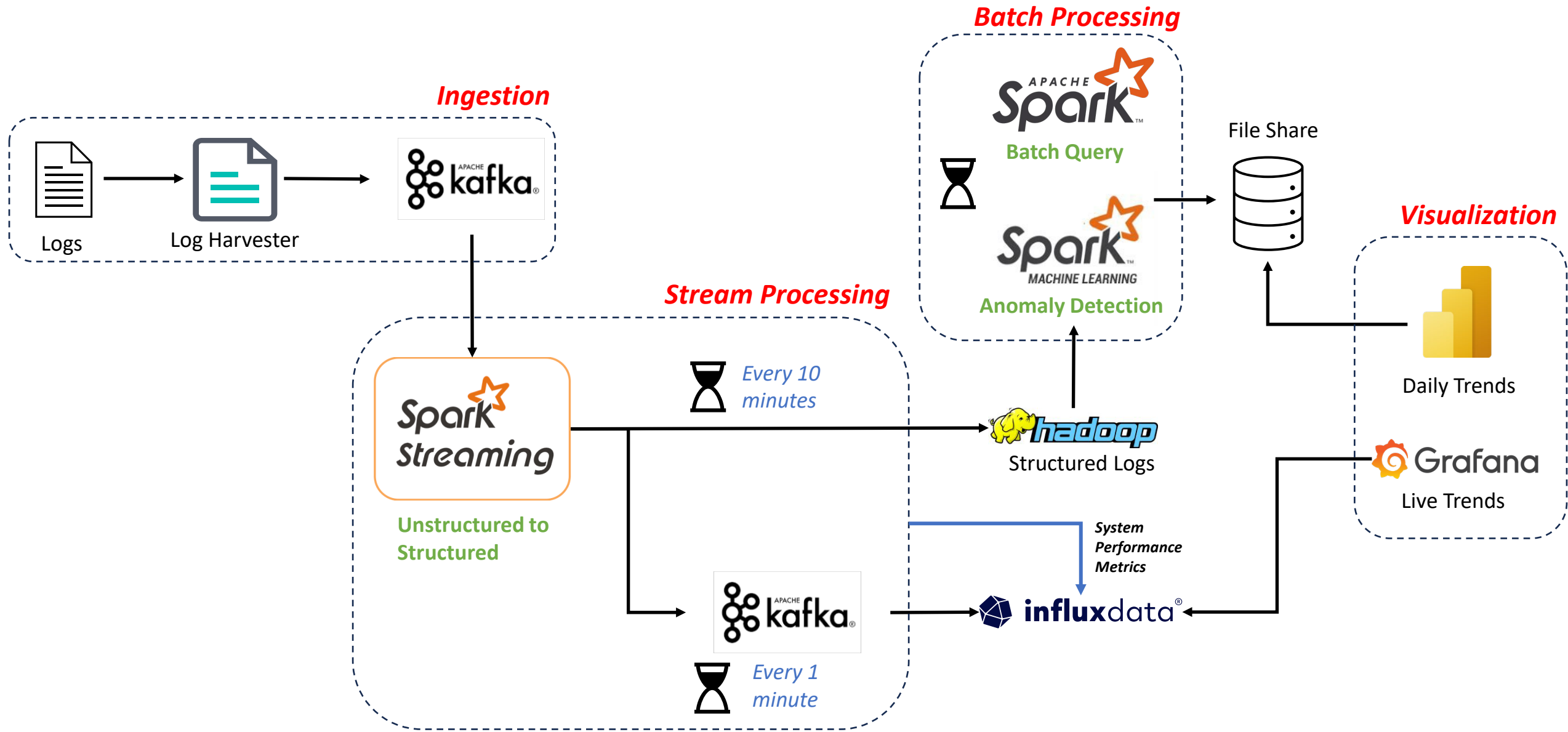
- **Distributed Data Parallel Processing**

- Efficiently handle large-scale log data by splitting tasks across multiple nodes
- Enable real-time monitoring and faster issue detection during high data volumes
- Enhance fault tolerance with redundancy and data replication to prevent data loss

- **Visualization of System KPIs**

- Present key metrics (failure rates) via interactive dashboards
- Track system performance and detect errors, malfunctions, or bottlenecks
- Visualize real-time data using charts
- Enable quicker anomaly detection and informed decision-making
- Improve operational monitoring, system stability, and performance optimization

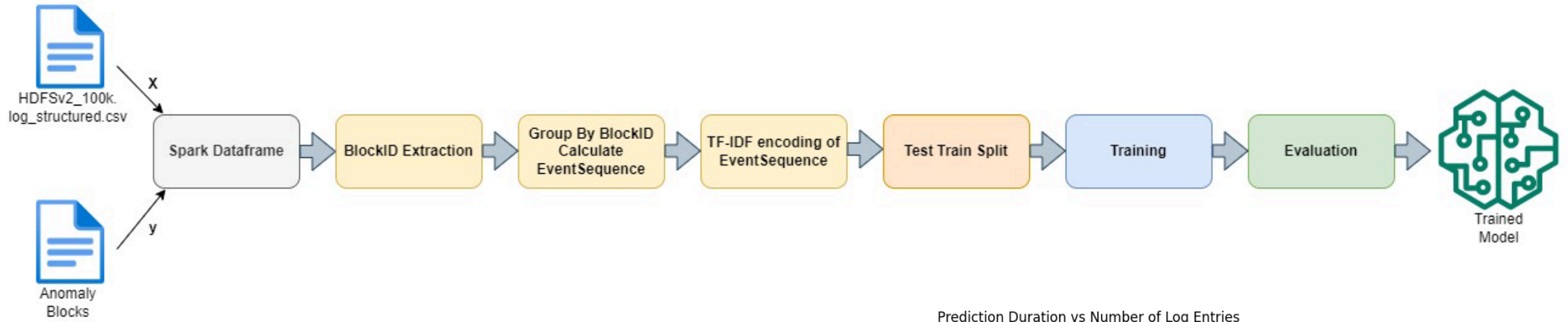
Architecture



Unstructured to Structured

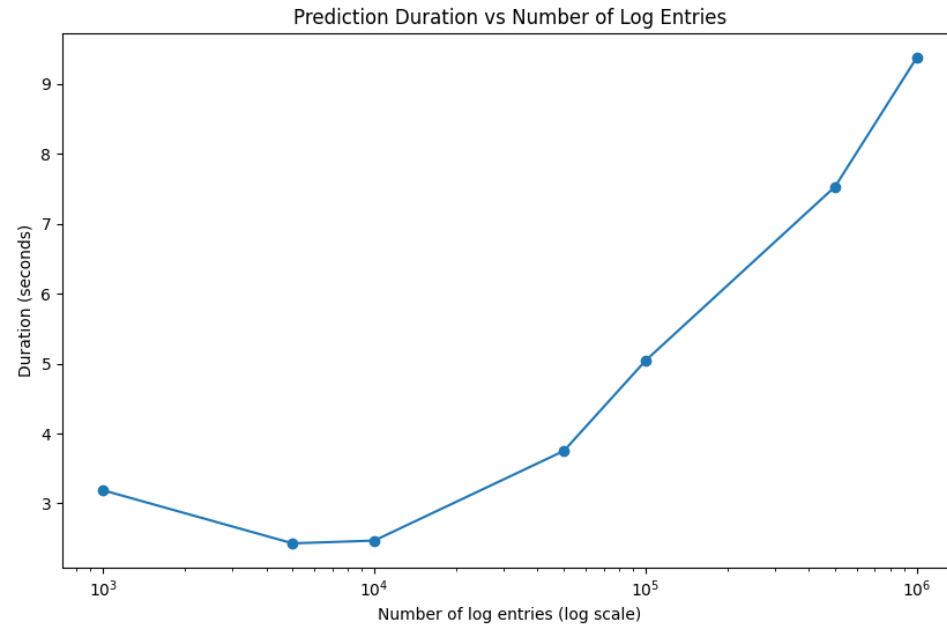


Anomaly Detection



Name	Value
Training duration	106.27 s
Accuracy	0.7723
Precision	0.9714
Recall	0.5151
F1 Score	0.6733
ROC-AUC	0.7723

Model Performance Metrics



Inference Performance Metrics

Development Setup



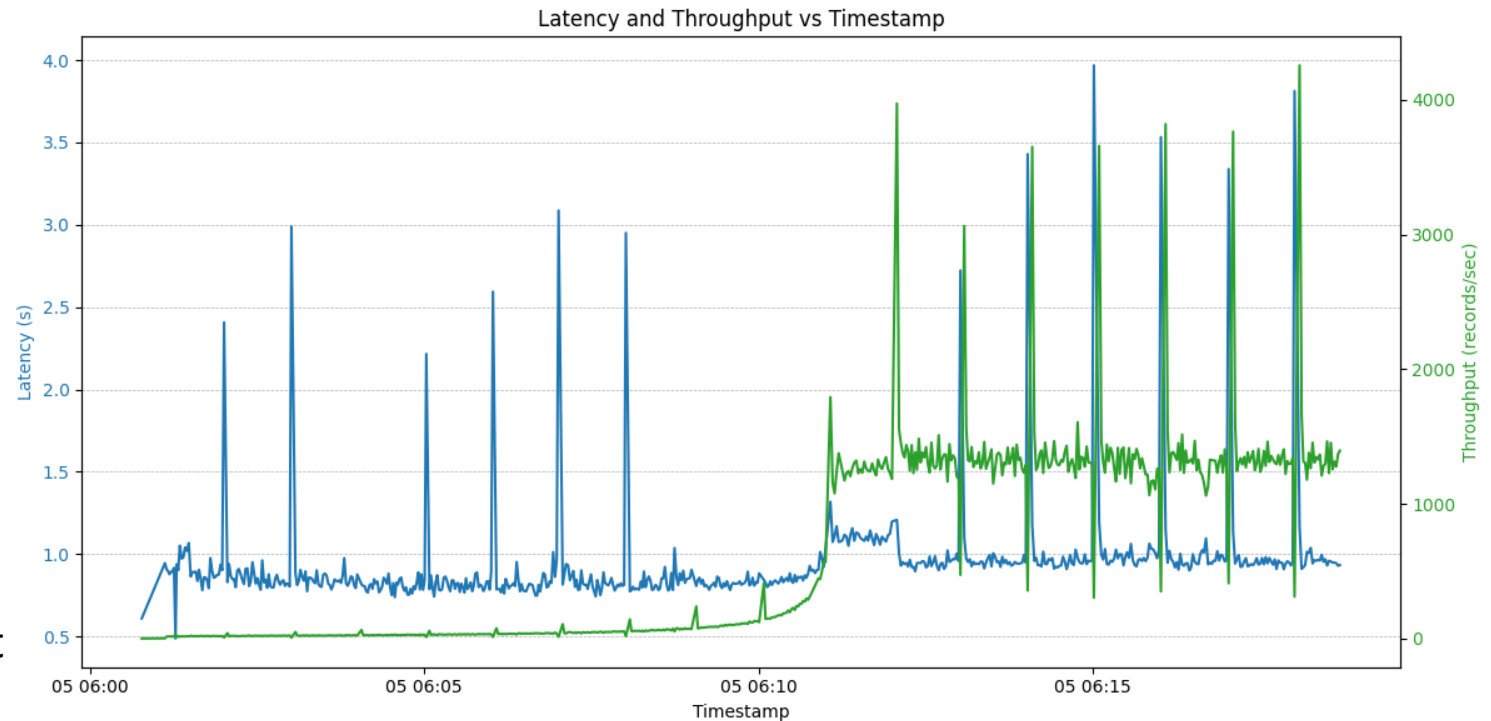
Evaluation: Throughput and latency

Experiment:

- Gradually increase input rate
- Measure latency of log parsing
- Measure throughput of the system

Observations:

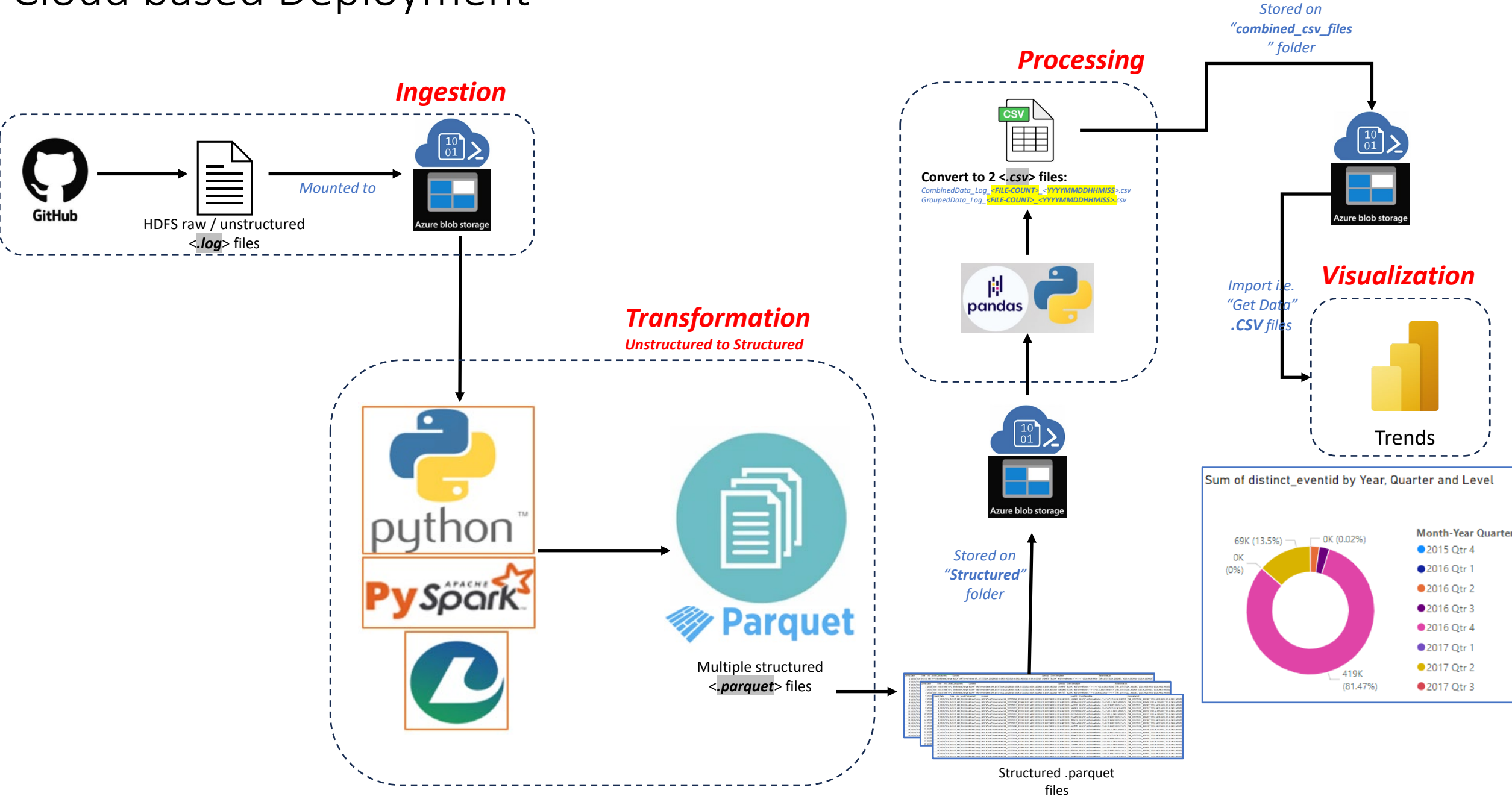
- No significant impact on latency
- Throughput saturates after a certain point



References and Further Reading

- 1. Jieming Zhu, Shilin He, Pinjia He, Jinyang Liu, Michael R. Lyu. [“Loghub: A Large Collection of System Log Datasets for AI-driven Log Analytics.”](#) (“GitHub - logpai/loghub: A large collection of system log datasets for ...”) (“GitHub - logpai/loghub: A large collection of system log datasets for ...”) IEEE International Symposium on Software Reliability Engineering (ISSRE), 2023.
- 2. Pinjia He, Jieming Zhu, Zibin Zheng, and Michael R. Lyu. (“LogAnMeta: Log Anomaly Detection Using Meta Learning - ResearchGate”) [“Drain: An Online Log Parsing Approach with Fixed Depth Tree](#), Proceedings of the 24th International Conference on Web Services (ICWS), 2017.” (“logparser/logparser/Drain/README.md at main - GitHub”)
- 3. Shilin He, Jieming Zhu, Pinjia He, Michael R. Lyu. [Experience Report: System Log Analysis for Anomaly Detection](#), *IEEE International Symposium on Software Reliability Engineering (ISSRE)*, 2016. (“GitHub - logpai/loghub: A large collection of system log datasets for ...”) [中文版本](#) *ISSRE Most Influential Paper*
- 4. [Use open source Drain3 log-template mining project to monitor for network outages](#)
- 5. [Fingerprinting the Datacenter: Automated Classification of Performance Crises](#), by Peter Bodík, Moises Goldszmidt, Armando Fox, Hans Andersen. *Microsoft*
- 6. [loghub](#)
- 7. [Drain3](#)
- 8. [loglizer](#)

Cloud based Deployment



MS-Azure (DFS) Architecture

Architecture Components

- **Log Storage:**
 - Use Azure Blob Storage to store raw HDFS logs (.log files).
 - Logs are ingested periodically from on-premises or other cloud locations.
- **Log Parsing:**
 - Drain Parser (integrated in Azure Databricks) is used to parse unstructured logs into structured log templates and extract metadata.
 - Parsed output is saved as .parquet files in Blob Storage.
- **Data Processing:** *Azure Databricks (Spark/PySpark)*
 - Aggregates and transforms the .parquet files.
 - Generates summarized .csv outputs for visualization.
 - Performs batch-based or stream-based log processing.
- **Analytics and Visualization:**
 - Azure Blob Storage serves as the data source for Power BI. Power BI visualizes aggregated data (e.g., event trends, latency, error distribution).
- **Workflow Automation:**
 - Azure Data Factory or Azure Logic Apps is used for orchestrating data ingestion, processing, and export pipelines.

Reasons for Failure

- **Resource Constraints:**
 - Free-tier clusters have limited memory and compute resources. Parsing an 800 MB file and generating 2000 small .parquet files can quickly exhaust the available memory or disk I/O capacity.
- **Excessive File Partitioning:**
 - This can overwhelm the cluster's metadata handling and slow down processing due to excessive partitioning.
- **Drain Parser Complexity:**
 - Drain parsers often involve recursive algorithms, which can be computationally expensive. Processing large files can cause memory overflows in constrained environments.
- **Driver Overload:**
 - If the driver node processes excessive data or metadata, it may crash due to insufficient memory or compute power.

Explain: Drain Parser

The Drain Parser is a tree-based log parsing algorithm designed to parse unstructured log messages into structured formats efficiently. It works by identifying patterns in log messages and extracting structured templates and parameters.

Key Features:

- **Tree-Based Approach:** Drain builds a hierarchical tree with nodes that represent log message patterns. Each path in the tree corresponds to a structured log template.
- **Parameters Extraction:** It separates dynamic parameters (e.g., timestamps, IP addresses, etc.) from static text in log messages.
- **Similarity Threshold (st):** Controls how similar a new log message must be to an existing template to fit into the same structure.
- **Depth (depth):** Limits the depth of the tree for log grouping.

Workflow of Drain Parser:

- **Input:** Raw unstructured log messages.
- **Log Template Extraction:** Static parts of the logs are grouped into templates, while dynamic parts are parameterized.

Code Explanation

The code is a pipeline that processes unstructured log files, parses them into a structured format using the Drain log parser, and then saves the results as a structured .parquet file for further analysis.

The provided code:

- *Reads Unstructured Logs:*
 - Reads log files stored in the DBFS mount (...:/mnt/...).
 - Uses Spark to parallelize log processing.

- **Parses Logs Using Drain Parser:**
 - Each partition of the RDD is processed by the `process_file_train` function.
 - Drain Parser takes the log format, regular expressions (regex), and similarity threshold to parse the logs.
- **Processes Multiple Files:**
 - Logs from multiple files are processed and combined into a unified DataFrame (*aggregated_df*).
- **Writes Output to Parquet:**
 - The final structured DataFrame is written as a *.parquet* file to .../mnt/hdfs2logs/structured.
- **The structured format includes:**

Column Name	Data-Type	Description
Date	String	The date extracted from the log message.
Time	String	The time extracted from the log message.
ms	Integer	Milliseconds extracted from the log message.
Level	String	Log severity level (e.g., INFO, ERROR, WARN).
Component	String	Log component/module generating the message.
Content	String	The dynamic content of the log message.
EventTemplate	String	Template representing the static structure of the log message.
EventId	Integer	A unique identifier assigned to each log template by Drain Parser.