
Building a big data pipeline on AWS

Introduction to Big data

- Real time streaming data being captured at regular intervals of time from millions of IOT devices like sensors,clickstreams,logs from the device APIs and historical data from SQL databases.
 - To store the huge volumes of data with high velocity and veracity ,we need an efficient scalable storage system which is distributed across different nodes either in local or in cloud.
 - Here comes the Hadoop concept which can be classified into two groups -Storage and processing
 - Storage will be done in HDFS and processing is done using Map reduce
 - Initially the processing of this data used to be done in Java by writing mapreduce programs with controller services. Now a days lot of tools have been developed on top of map reduce such as Hive ,Presto,Druid etc which can be directly used by their built in SQL interfaces and are much faster due to inbuilt indexing of the data.
 - Underlying storage will be HDFS for all kinds of big data scenarios and underlying processing is Map reduce in most of the cases.
-

Introduction to Data Pipeline

- It refers to a system for moving **data** from one system to another. The **data** may or may not be transformed, and it may be processed in real time (or streaming) instead of batches.
 - Right from extracting or capturing data using various tools , storing raw data, cleaning, validating data, transforming data into query worthy format , visualisation of KPIs including Orchestration of the above process is data pipeline.
 - Scalability, reliability, security should be taken into consideration in each step of data pipeline building in big data environments.
 - Transformed data can be used to derive KPIs and derive initial insights from the data for exploratory analysis
 - Further analysis of data using ML algorithms can be used to derive predictions out of the data either historical or near real time .
-

Business Impact of building data pipelines

- Saves time and effort for actuaries, BI analysts
 - Manual process of cleaning, validating of raw data takes on an average of 3-4 hours of time for each and every raw file received at a certain interval of time.
 - Automating this simple process saves 3-4 hours for each run.
 - Similarly automation of ETLs using Cron, Oozie, Airflow processes etc saves effort and time for Data engineers as well as Analysts.
 - Only pipeline monitoring and maintenance needs to be taken care of once the pipeline is built.
 - Grabbing the data from source to destination-sources being APIs, Google sheets, Local systems, SQL based databases etc. Destination systems can be HDFS for storage, Kafka, Flume, Storm, Logstash to grab real time streaming data
 -
-

Few source and destination data systems

Source systems

- Streaming APIs
- SQL databases
- Google sheets
- Local systems

Data capture

Kafka
Flume
Logstash
NiFi
Storme
Python

Destination systems

HDFS followed by Hive
Processing using Pyspark
HiveQL
Elasticsearch
Druid SQL

Visualization tools

Tableau
PoweBI
AWS Quicksight
Kibana
Grafana
Qlikview
Python

Overview

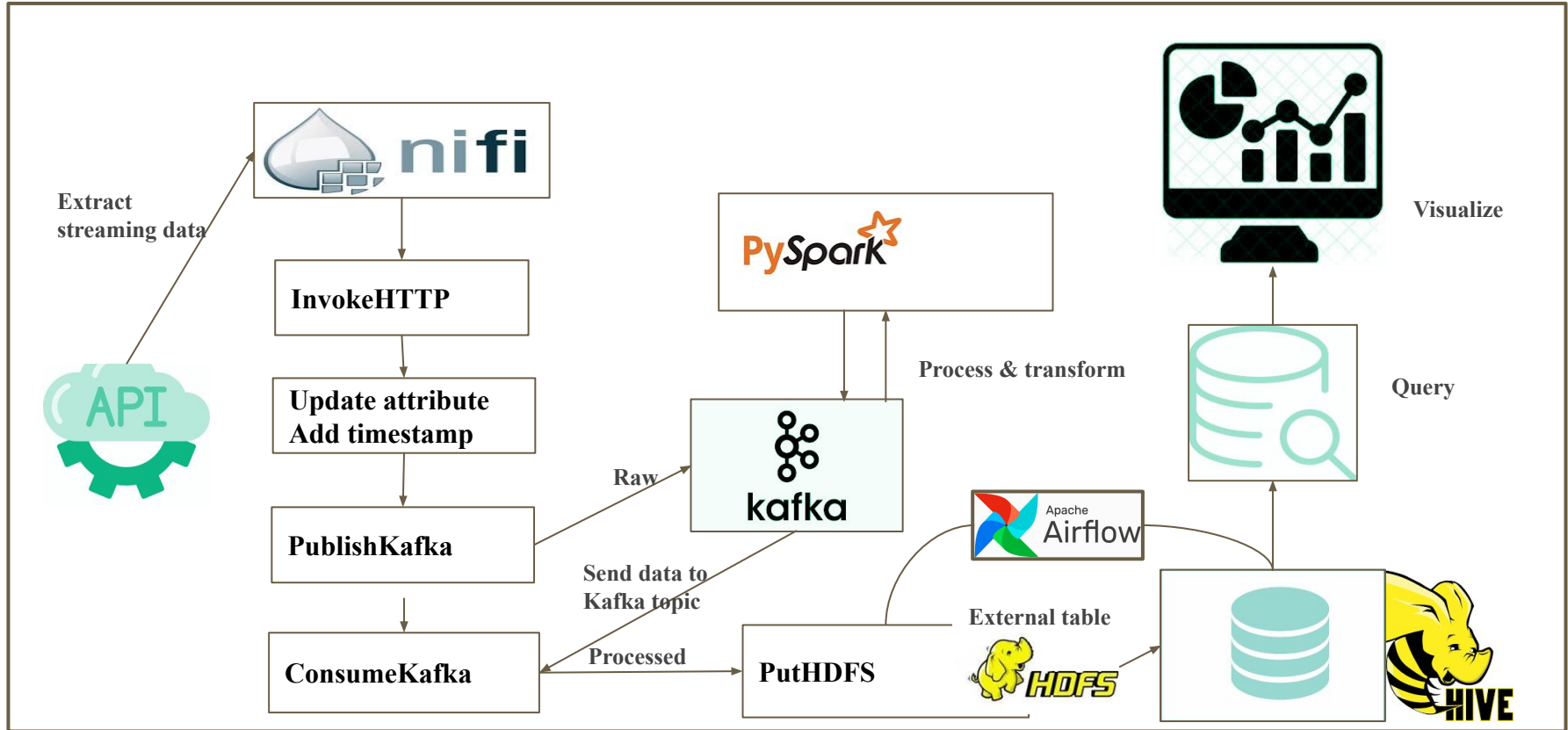
Process to build a data pipeline to automate data ingestion & processing

- Extract streaming data from APIs at a certain interval of time using NiFi.
- Parse the data using Nifi which pushes the data into Kafka
- Extract the data from Kafka into PySpark and process it .
- Write the streaming output data into another Kafka topic using NiFi. Put the processed streaming data into HDFS.
- Create external table in Hive /Presto (Druid is more suggestible).
- Query the data from Hive/Presto/Druid. (Perform more ELTs)
- Orchestrate the process.

System requirements

- Ec2 instance in AWS - t2.xLarge 16GB of RAM and Ubuntu 16.04 image on it , installation steps attached:
- <https://docs.google.com/document/d/1sot7fCIgLZmldpcjodbQ0RZPRwN367-aOpa9d9UJS5M/edit?usp=sharing>
- Good internet connectivity.
- Install required services required for Data architecture on Ubuntu machine installed above , installation steps attached:
- <https://docs.google.com/document/d/1UcN0ycXVPuOosGdMxz5g6WLgqneeY4z-vWAz15zknG4/edit?usp=sharing>
- Services included in documentation are :
- Hadoop
- Hive
- Scala- Spark
- Zookeeper
- Kafka
- NiFi
- Airflow

Data architecture



Apache NiFi

- Open source for automating and managing data flow between systems.
- Process distributed data and executes on the JVM on host OS.
- Web based UI for creating, monitoring & controlling data flows.
- Supports buffering of all Queued data.
- Processors connect to many source and destination systems.
- Easy to troubleshoot and flow optimisation.
- Role based authentication
- Build user defined processors and controller services.
- Easy of use- need not code much.



Apache Kafka

- Apache Kafka is an open-source stream-processing software which buffers the records as they occur without data loss
- Uses I/O efficiently by batching and compressing records
- Supports streaming data and is scalable, durable and fault tolerant
- Uses file system for caching and storage
- Works on publish and subscribe mechanism
- Producers are source systems , consumers are destination systems
- Messages stored in topic
- Why Kafka over Apache Flume or Apache Storm?



Kafka Vs Flume

Kafka	Flume
Used as a bridge between source and destination systems	Mostly used to gather data into Hadoop from various sources
No data loss and more general purpose system can be used on any data	No data loss but mostly used for log based aggregations and log data
Multiple publishers and subscribers can share same topic	It's a tool to gather data into HDFS
Makes data available even after single point failure	Depends on the configuration

Challenges with Hive

- Long time to aggregate data in the warehouse using Hive or Presto at Query time.
- Once the velocity of the data increases , infrastructure to hold the data should be scalable and increased which results in more costs.
- But can be useful for complex calculations -So many architectures involve Hive for data computation and Presto/Druid for Querying the computed data.

Hive Vs Presto Vs Druid

Hive	Presto	Druid
Data warehousing package for Hadoop	Distributed SQL engine which doesn't store or hold any data	Druid is a column-oriented, open-source, distributed data store written in Java
Optimized for query throughput	Optimized for query latency	Optimized for very high query latency
All computations are performed as map reduce operations which takes more time	Computations are performed in-memory and doesn't use map reduce	Architecture of separating data storage and caching temporarily makes Druid very fast querying engine.
Aggregations cannot be performed during data ingestion.	Aggregations are performed only during query time -but much faster than Hive	Ingest time aggregations are supported. Hence do not take much time to query aggregated data.

Apache Spark

- **Apache Spark** is a unified analytics engine for big data processing, with built-in modules for streaming, SQL, machine learning and graph processing.
- An ETL engine which can run multiple distributed workloads.
- Can integrate with various databases like Redshift, Druid, SQL based databases -has connectors to them
- Has Python interface known as PySpark
- Provides SQL based interface known as Spark SQL
- Execution is based on DAGs(Directed Acyclic Graphs)
- Tasks are executed in memory
- 20X faster than Hive



Apache Airflow Vs Oozie

- Apache Airflow is an open-source workflow management platform. It started at Airbnb in October 2014 as a solution to manage the company's increasing complex workflows.
- Creating Airflow allowed Airbnb to programmatically author and schedule their workflows and monitor them via the built-in Airflow user interface
- Mainly used for data orchestration (automation)
- Provides better interface than Oozie
- Helps to build data pipelines in real time and supports various tasks like Hive ,Python ,SQL,email tasks etc.
- An **Airflow** workflow is designed as a directed acyclic graph (DAG). That means, that when authoring a workflow, you should think how it could be divided into tasks which **can** be executed independently.
- **Airflow** leverages growing use of python to allow you to create extremely complex workflows, while **Oozie** allows you to write your workflows in Java and XML. The open-source community supporting **Airflow** is 20x the size of the community supporting **Oozie**.

