

Introduction to Streaming & Messaging

Flume ,Kafka,SQS,
Kinesis streams & firehose
Kinesis Analytics

Omid Vahdaty, Big Data Ninja



What is batch Processing?

the execution of a series of programs each on a set or "**batch**" of inputs, rather than a single input (which would instead be a custom job)



What is Streaming ?

Streaming Data is data that is generated continuously by thousands of data sources, which typically send in the data records simultaneously, and in small sizes (order of Kilobytes)



Streaming VS. Batch Processing

	Batch	Stream
Data Scope	Query the entire batch, with slight delay	Query most recent events defined in a time window.
Data Size	Large data sets	A few Individual records
Latency?	Minutes ,hours	Seconds, Milliseconds
Analysis	Complex Analytics	Basic: aggregations, metrics etc.



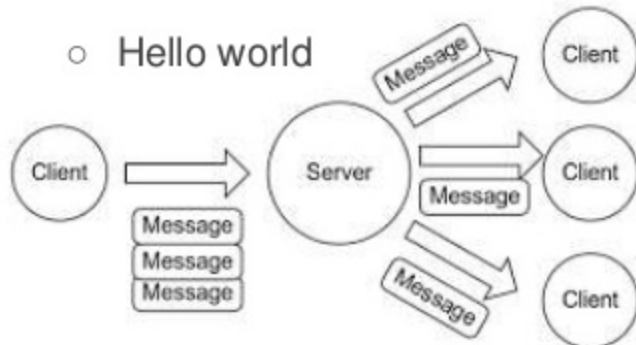
Challenges with Streaming Data

- Processing layer
 - Consuming data
 - Processing data
 - Notifying storage layer what to do.
- Storage layer
 - Ordering mechanism
 - Strong Consistency mechanism
- In general MUST have features:
 - scalability
 - data durability
 - fault tolerance



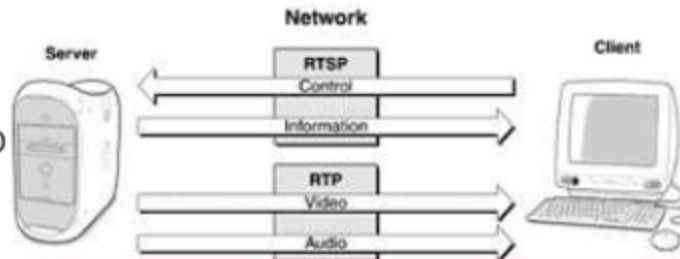
Messaging VS Streaming?

- Messaging: **framed message** based protocol.
- e.g 3 messages sent will look like:
 - Hello world
 - Hello world
 - Hello world



- Streaming: **unframed data** (bytes) stream based protocol
- e.g 3 messages sent will look like:

- Hell
- ow wo
- rld Hel
- low wor
- ldHellow wo
- rld



Messaging

Open Source: Kafka, flume
AWS: SQS

Flume

Flume Pros:

Good documentation with many existing implementation patterns to follow

Easy integration with existing monitoring framework

Integration with Cloudera Manager to monitor Flume processes

Flume Cons:

Event rather than stream centric

Calculating capacity is not an exact science but rather confirmed through trials

Throughput is dependent on the channel backing store.

Flume lacks the clear scaling and resiliency configurations (trivial with Kafka and Kinesis)



Kafka

Kafka Pros:

High achievable ingest rates with clear scaling pattern

High resiliency via distributed replicas with little impact on throughput

Kafka Cons:

No current framework for monitoring and configuring producers



Flume VS. Kafka

	Flume	Kafka
Choose when you desire	No need for customization. Need out of the box components such HDFS sink	Need a custom made high availability delivery system
Velocity	high	higher
Event processing		

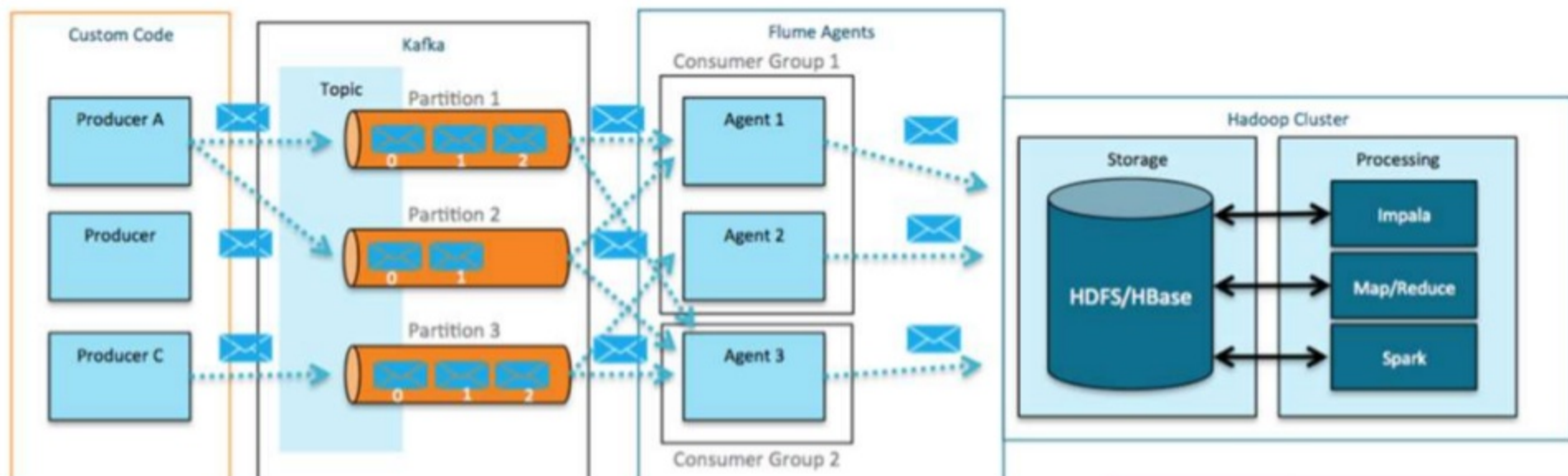


	Flume	Kafka
Original Motivation	distributed, reliable, and available system for efficiently collecting, aggregating and moving large amounts of log data from many different sources to a centralized data store. Built around hadoop ecosystem	general purpose distributed publish-subscribe messaging system Multi-consumer ultra-high availability messaging system.
Data Flow	push	pull
event availability	JDBC Databases Channel, file Channel. Loose flume agent = losing data.	replication of your events data by design.
Commercial support	Cloudera	Cloudera
Collectors built in	Yes.	just the messaging

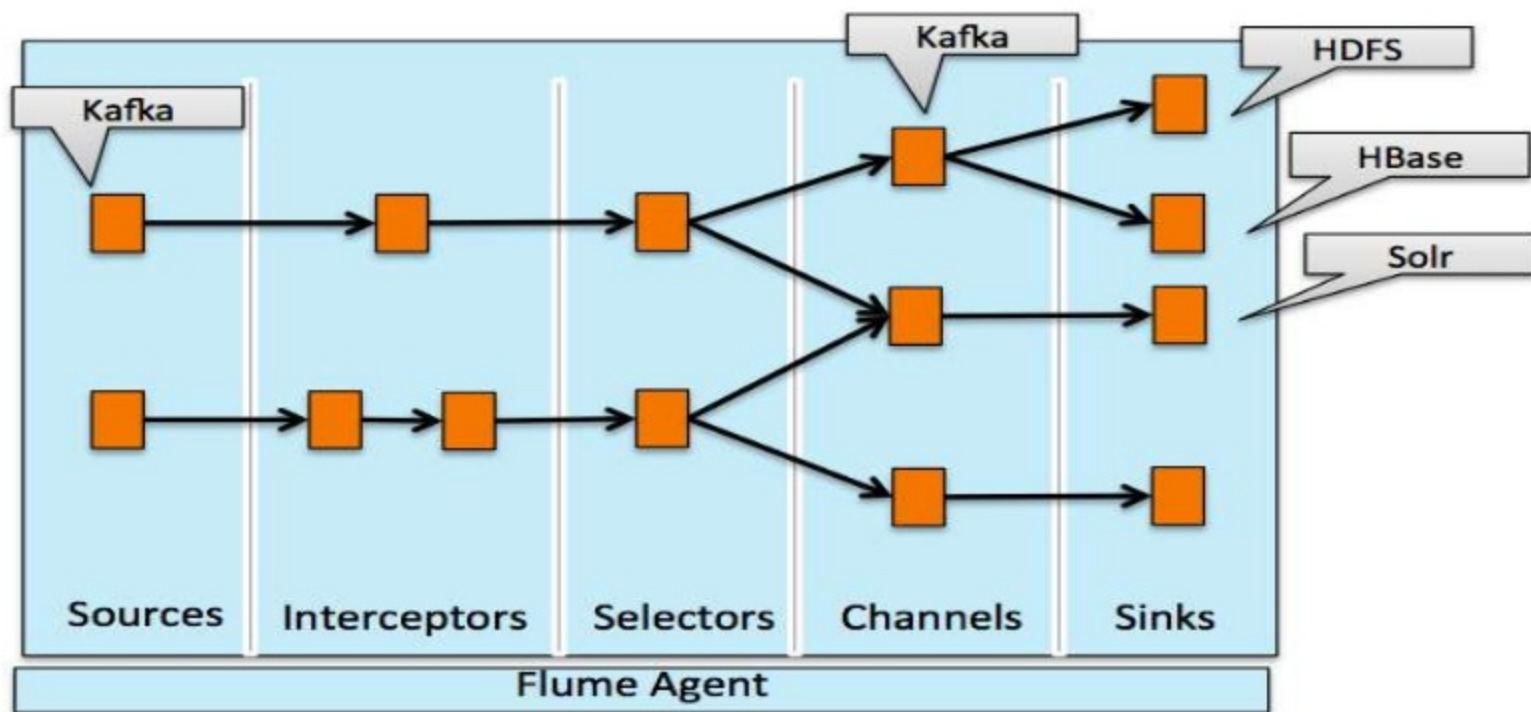


Use Case: Kafka and Flume combined

- Flume supports: Kafka source, Kafka channel, Kafka sink
- So, take the advantage of both and combine them to your needs.



Use Case: Kafka as a Channel



AWS SQS

- a fast, reliable, scalable, fully managed **message queuing** service
- **decouple the components** of a cloud application, move data between diverse, distributed application components without losing messages and without requiring each component to be always available.
- high throughput and **at-least-once processing**, and **FIFO queues**
- all messages are stored **redundantly** across multiple servers and data centers.
- Start with **three API calls : SendMessage, ReceiveMessage, and DeleteMessage**. Additional APIs are available to provide advanced functionality.
- **Queues**
 - **Standard queues** offer **maximum throughput**, best-effort ordering, and at-least-once delivery.
 - **FIFO queues** are designed to ensure **strict ordering** and exactly-once processing, with limited throughput.
- scales dynamically
- Authentication mechanisms


AWS SQS use cases

Messaging semantics (such as **message-level ack/fail**) and visibility timeout. For example, you have a queue of work items and want to track the successful completion of each item independently. Amazon SQS tracks the ack/fail, so the application does not have to maintain a persistent checkpoint/cursor. Amazon SQS will delete acked messages and redeliver failed messages after a configured visibility timeout.

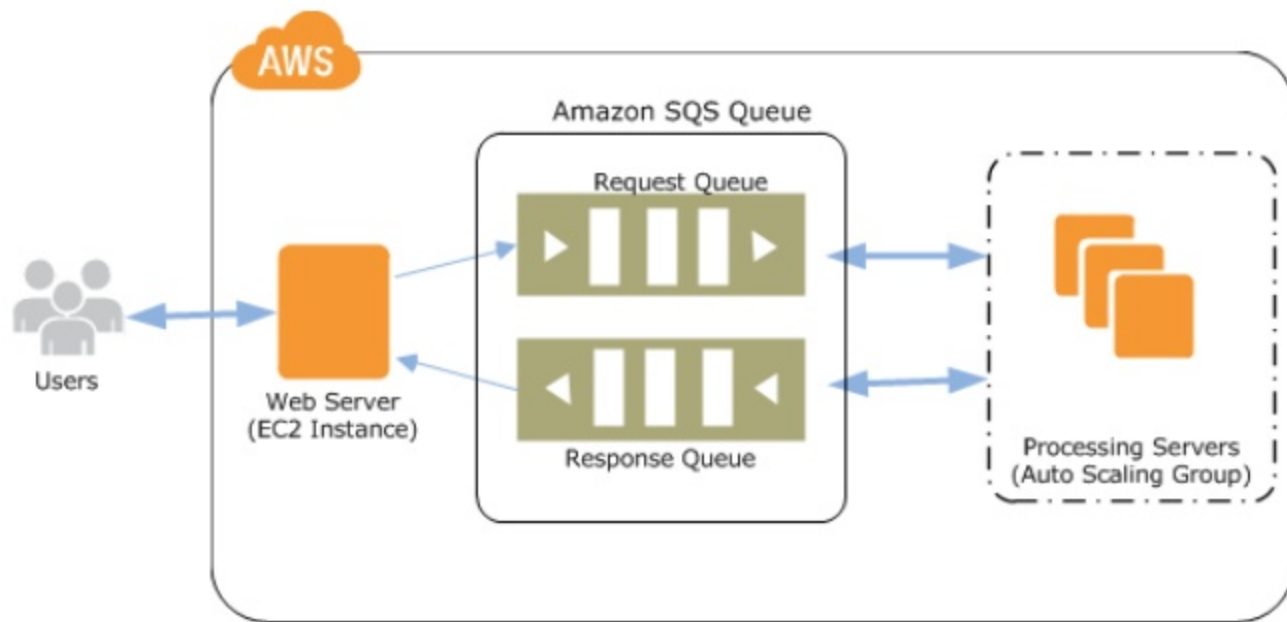
Individual message delay. For example, you have a job queue and need to schedule individual jobs with a delay. With Amazon SQS, you can configure individual messages to have a delay of up to 15 minutes.

Dynamically increasing concurrency/throughput at read time. For example, you have a work queue and want to add more readers until the backlog is cleared. With Amazon Kinesis, you can scale up to a sufficient number of shards (note, however, that you'll need to provision enough shards ahead of time).

Leveraging Amazon SQS's ability to scale transparently. For example, you buffer requests and the load changes as a result of occasional load spikes or the natural growth of your business. Because each buffered request can be processed independently.



Typical SQL use case: decoupling APP layers.



Streaming

AWS Kinesis:
Streams, Firehose, Analytics



AWS Kinesis (streams)

- build custom applications that process or analyze streams
- continuously capture and store terabytes of data per hour
- Hundreds sources
- allows for real-time data processing
- Easy to use, get started in minutes
 - Kinesis Client Library
 - Kinesis Producer Library
- allows you to have multiple Applications processing the same stream concurrently.
- The throughput can scale from megabytes to terabytes per hour
- synchronously replicates your streaming data across three AZ
- preserves your data for up to 7 days



AWS Kinesis (streams) use cases

- Log and Event collection
- Mobile Data collection
- Real Time Analytics
 - when loading data from transactional databases into data warehouses.
 - Multi-stage processing using specialized algorithms
 - stream partitioning for finer control over scaling
- Gaming Data feed



AWS Kinesis (streams) use cases

Routing related records to the same record processor (as in streaming MapReduce). For example, **counting and aggregation are simpler when all records for a given key are routed to the same record processor.**

Ordering of records. For example, you want to **transfer log data** from the application host to the processing/archival host while maintaining the order of log statements.

Ability for multiple applications to **consume the same stream concurrently**. For example, you have one application that updates a real-time dashboard and another that archives data to Amazon Redshift. You want both applications to consume data from the same stream concurrently and independently.

Ability to **consume records in the same order a few hours later**. For example, you have a billing application and an audit application that runs a few hours behind the billing application. Because Amazon Kinesis stores data for up to 24 hours, you can run the audit application up to 24 hours behind the billing application.



AWS Kinesis (streams)

Kinesis Pros:

High achievable ingest rates with clear scaling pattern

Similar throughput and resiliency characteristics to Kafka

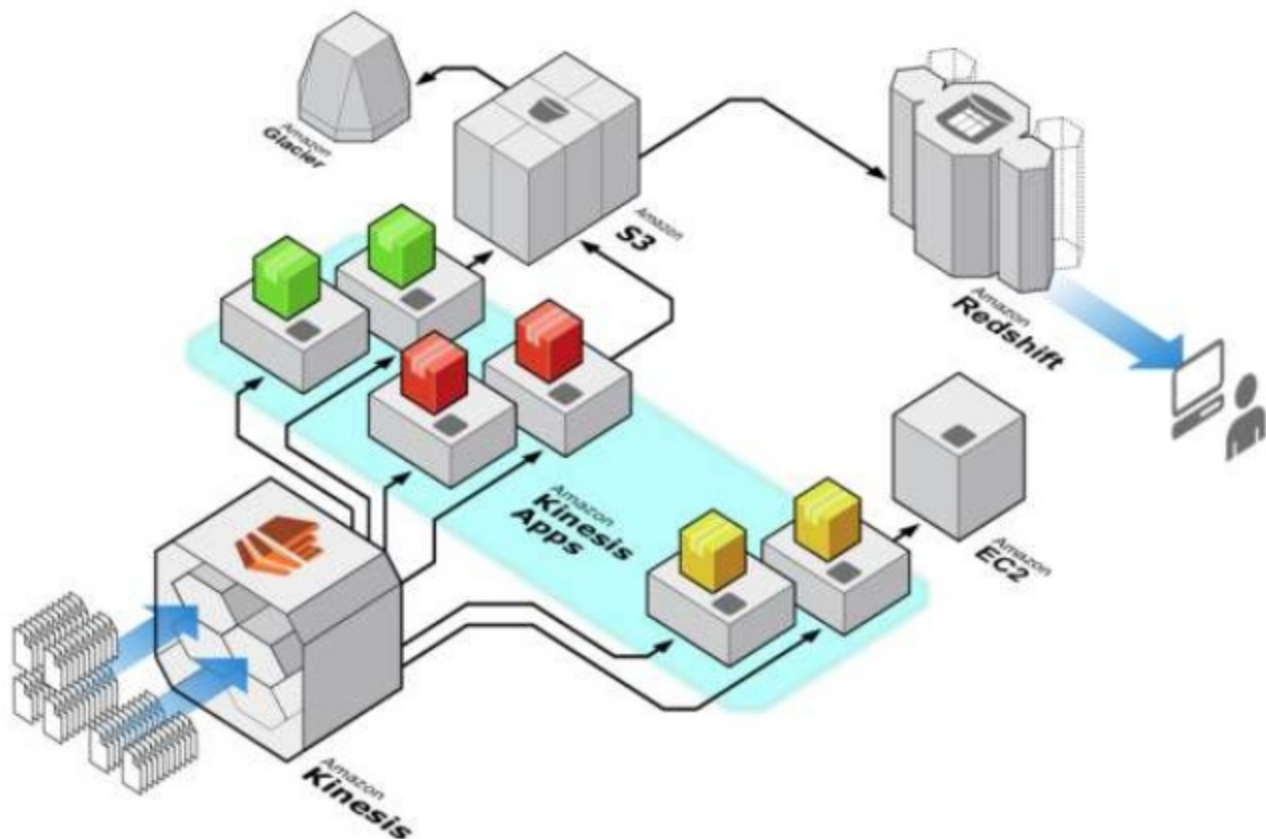
Integrates with other AWS services like EMR and Data Pipeline.

Kinesis Cons:

No current framework for monitoring and configuring producers

Cloud service only. Possible increase in latency from source to Kinesis.

AWS Kinesis (streams)



AWS Kinesis Firehose



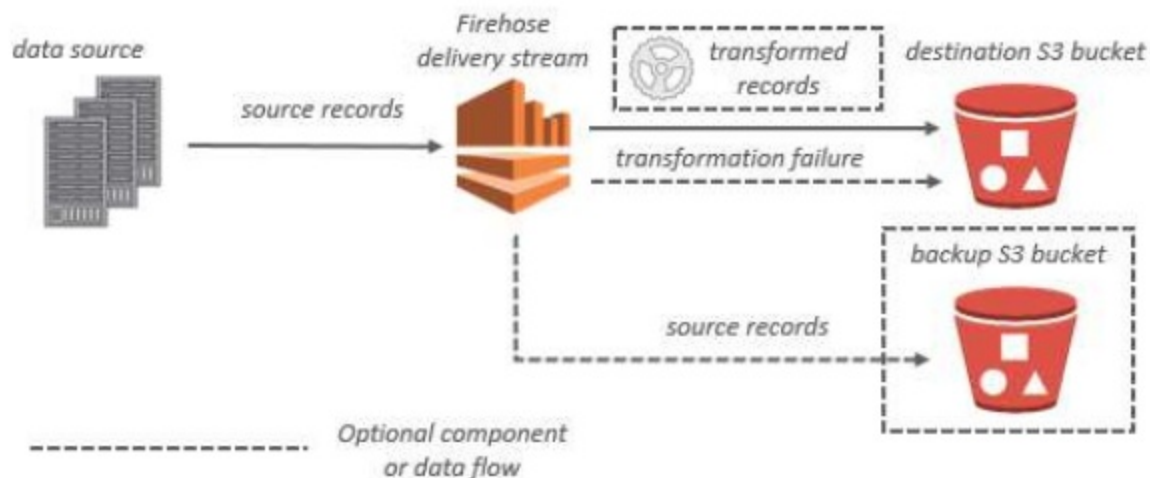
- the easiest way to load streaming data into AWS.
- **capture, transform, and load streaming data**
 - integrates into Kinesis Analytics, S3, Redshift, Elasticsearch Service
 - Serverless Transformation on RAW data. (lambda function)
 - E.g transform log file into CSV format
- Firehose can back up all untransformed records to your S3 bucket concurrently while delivering transformed records to the destination. You can enable source record backup
- enabling near real-time analytics
- Easy to use.
- Monitoring options.
- Limits
 - 20 stream per regions, Each stream
 - 2000 transaction per sec
 - 5000 records per sec
 - 5MB/s
 - Support 24 hours replay in cases on downtime

Kinesis Firehose agent

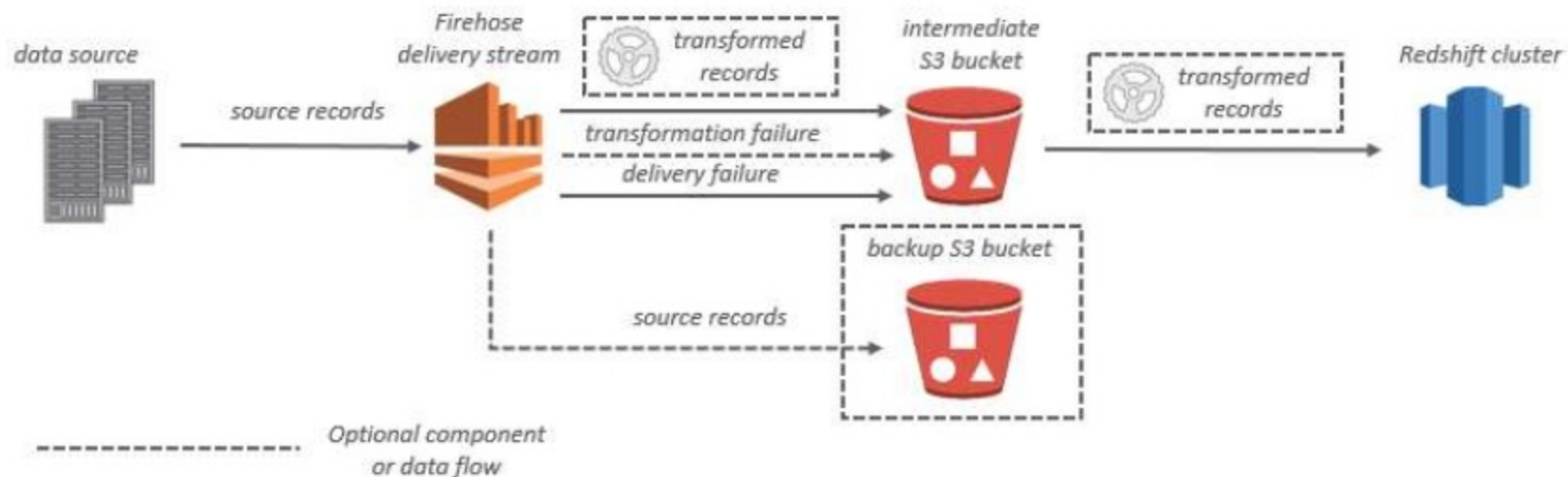
- Java software app that send data to streams/firehose
- monitors a set of files for new data and then sends streams/firehose
- It handles file rotation, checkpointing, and retrial upon failures.
- supports [Amazon CloudWatch](#) so that you can closely monitor and troubleshoot the data flow from the agent.
- Data processing options:
 - SINGLELINE – This option converts a multi-line record to a single line record by removing newline characters, and leading and trailing spaces.
 - CSVTOJSON – This option converts a record from delimiter separated format to JSON format.
 - LOGTOJSON – This option converts a record from several commonly used log formats to JSON format. Currently supported log formats are [Apache Common Log](#), Apache Combined Log, Apache Error Log, and [RFC3164](#) (syslog).
- <https://github.com/aws-labs/amazon-kinesis-agent>

Write a JAVA agent to Firehose

- AWS java SDK
- Firehose API
 - Single record: [PutRecord](#)
 - Batch: [PutRecordBatch](#).
- Key concepts:
 - **Firehose delivery stream**
 - **Data producer** - i.e web server creating log.
 - **Record**: The data of interest that your data producer sends to a Firehose delivery stream. A record can be as large as 1000 KB.
 - **buffer size** (in MB)
 - **buffer interval** (seconds)
- **Java examples**: <http://docs.aws.amazon.com/firehose/latest/dev/writing-with-sdk.html>



Firehose and Redshift usecase

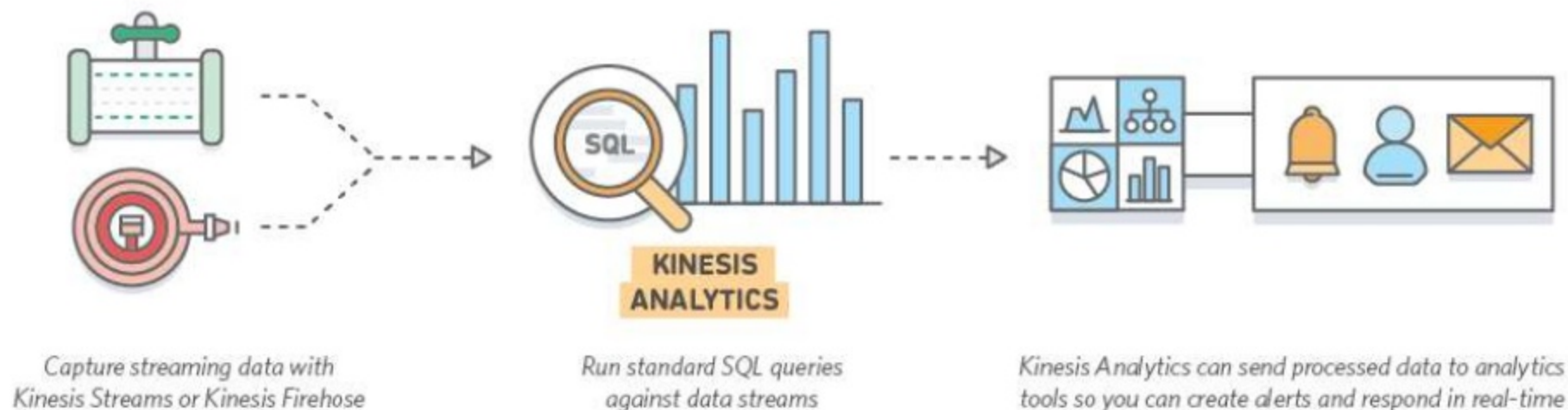


Kinesis Analytics : in-flight analytics.




- process streaming data in real time with standard SQL
- Amazon Kinesis Analytics enables you to create and run SQL queries on streaming c
- Easy 3 steps
 1. Configure Input stream (kinesis stream, kinesis firehose)
 - a. Automatically created Schema
 - b. Manually change schema if you like
 2. Write SQL query
 3. Configure output stream: s3, redshift, elastics search
- Elastic: scale up down
- Managed service
- Standard SQL

Kinesis Analytics : in-flight analytics.



Stay in touch...

- [Omid Vahdaty](#) 
- +972-54-2384178



 **HALO**
ANALYTICS



The logo for Jajah Telefonica. It features the word 'jajah' in a stylized, pink, sans-serif font, with a small, curved line above the 'j'. Below 'jajah' is the word 'Telefonica' in a black, cursive script font.