

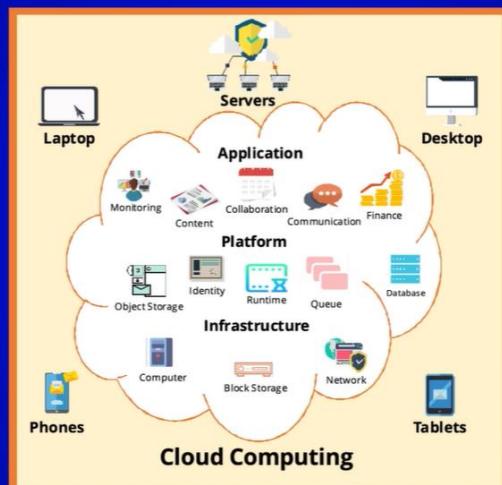
Learning Objectives

After completing this lesson, you should be able to understand:

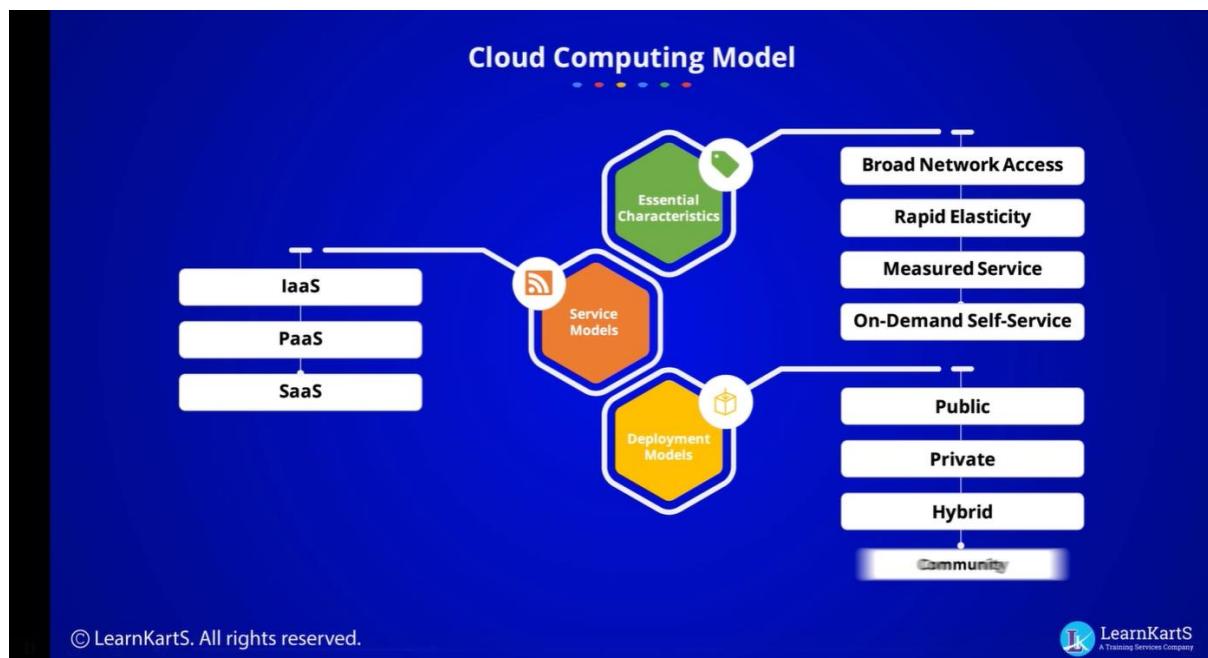
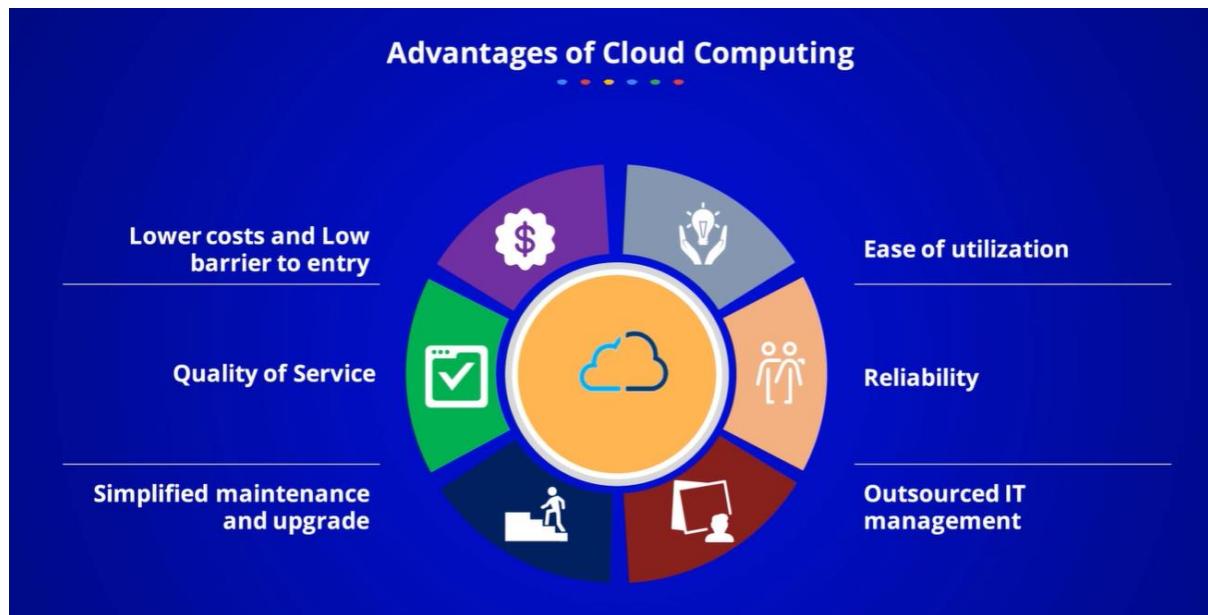


Cloud computing introduction, advantages and types.

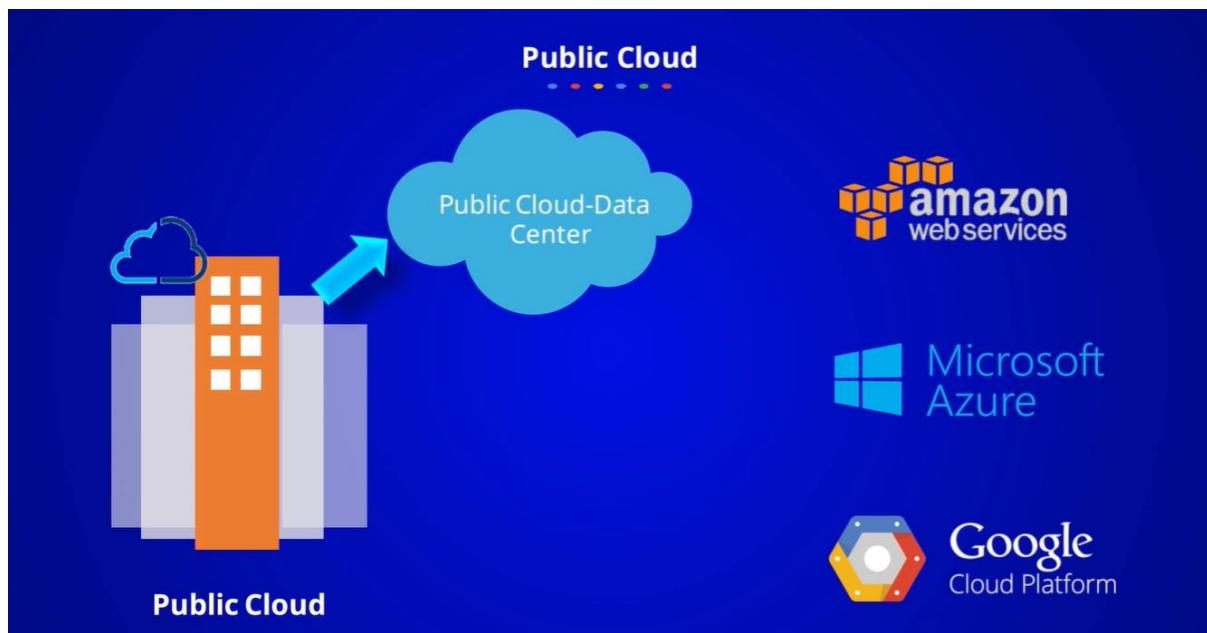
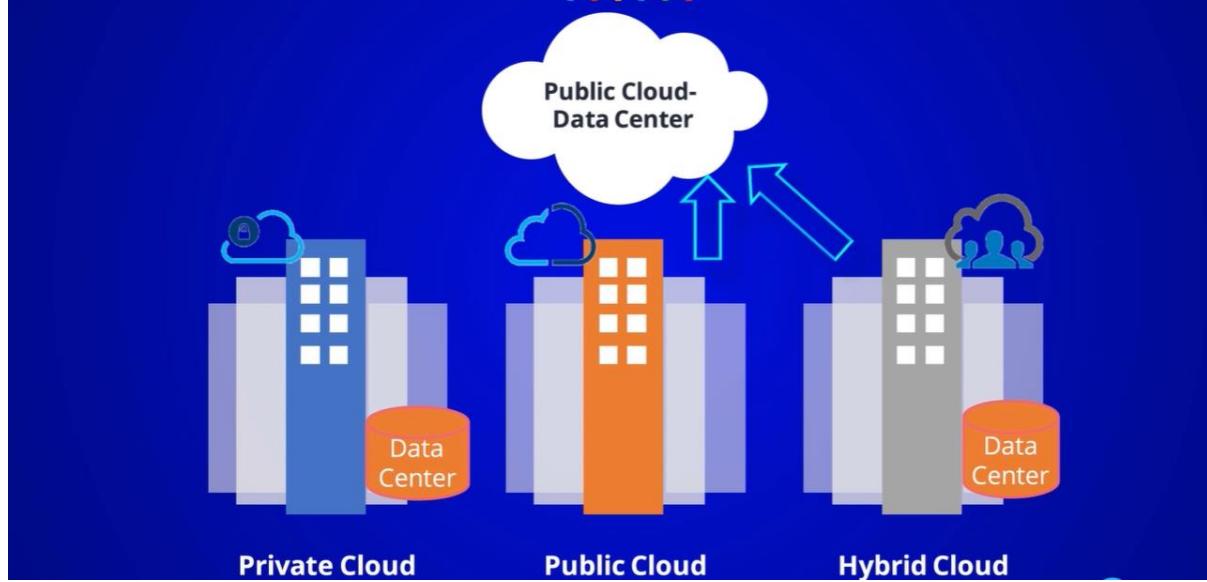
What is Cloud Computing?



Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a common pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be quickly provisioned and released with minimal management exertion or service provider interaction.



Cloud Computing Deployment Model (Cont..)



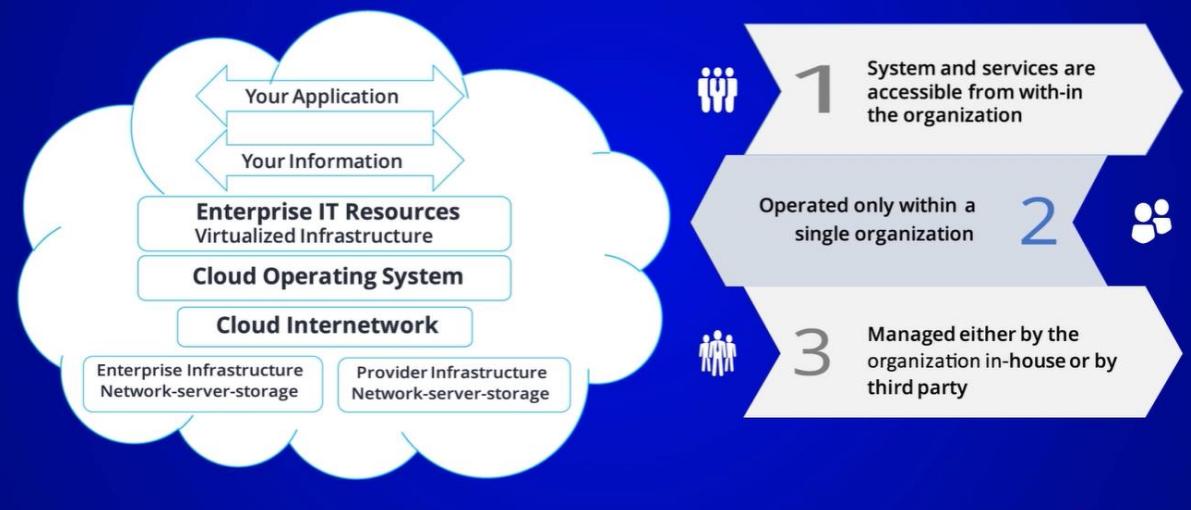
Public Cloud Advantages

Public Cloud Key Characteristics

-  Available to general public
-  Located at the premise of the cloud provider
-  May be owned by a private company, organization, academic institution or a combination of owners
-  Easy setup and inexpensive to the customer
-  Pay only for the services consumed

Private Cloud

Private Cloud Model



Private Cloud Disadvantages

Restricted Area of Operation

Private cloud is only accessible with-in the organization hence cannot be deployed globally.

Limited Scalability

Scalability is limited to the resources which are available with-in organization

High Priced

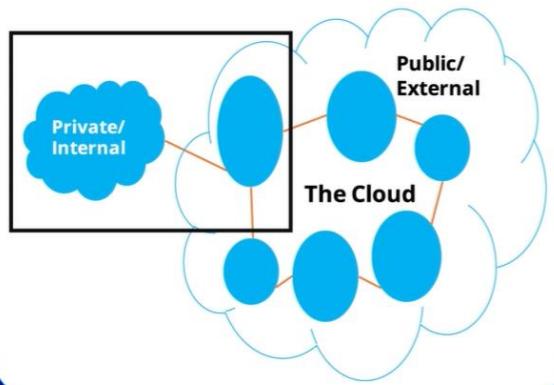
Organizations need to purchase hardware to deploy a Private cloud

Additional Skills

Maintaining cloud infrastructure required additional expertise with-in organization

Hybrid Cloud

Hybrid Cloud Model



Mixture of **public** and **private** cloud.

Non-critical activities are performed using public cloud



Critical activities are performed using private cloud



Hybrid Cloud Advantages



A stylized illustration of a person in a blue shirt and dark pants climbing a light blue mountain. Three circular icons are placed on the mountain's peak, each connected by a line to a text label: an orange icon with a square pattern is connected to 'Composed of two or more cloud models'; a teal icon with a laptop-like shape is connected to 'Portability'; and a green icon with three arrows forming a triangle is connected to 'Load balancing, high availability, Disaster Recovery'.

Composed of two or more cloud models

Portability

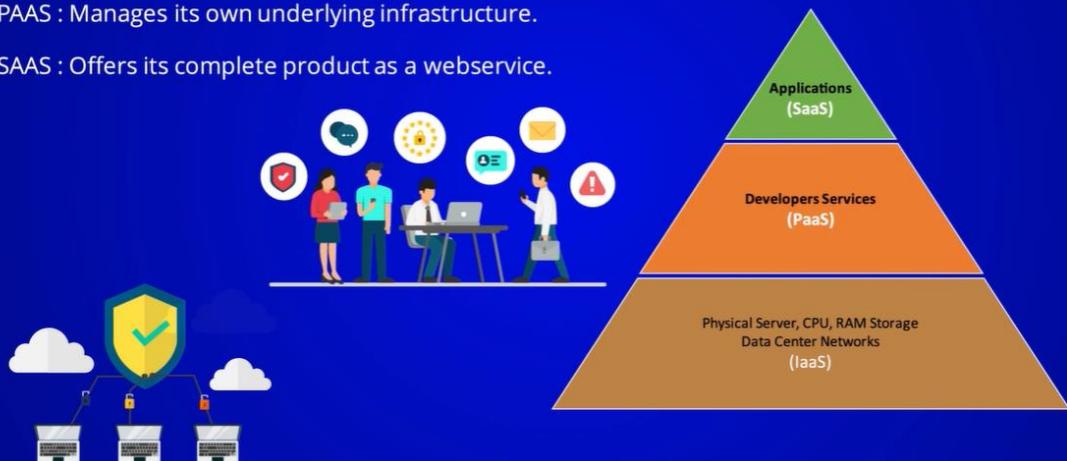
Load balancing, high availability, Disaster Recovery

Cloud Service Categories

IAAS : Provides Basic Building Block for Cloud IT

PAAS : Manages its own underlying infrastructure.

SAAS : Offers its complete product as a webservice.



Cloud Service Categories

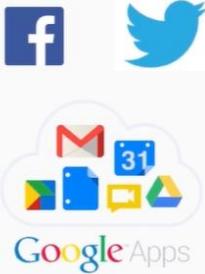
IaaS (Infrastructure as a Service)

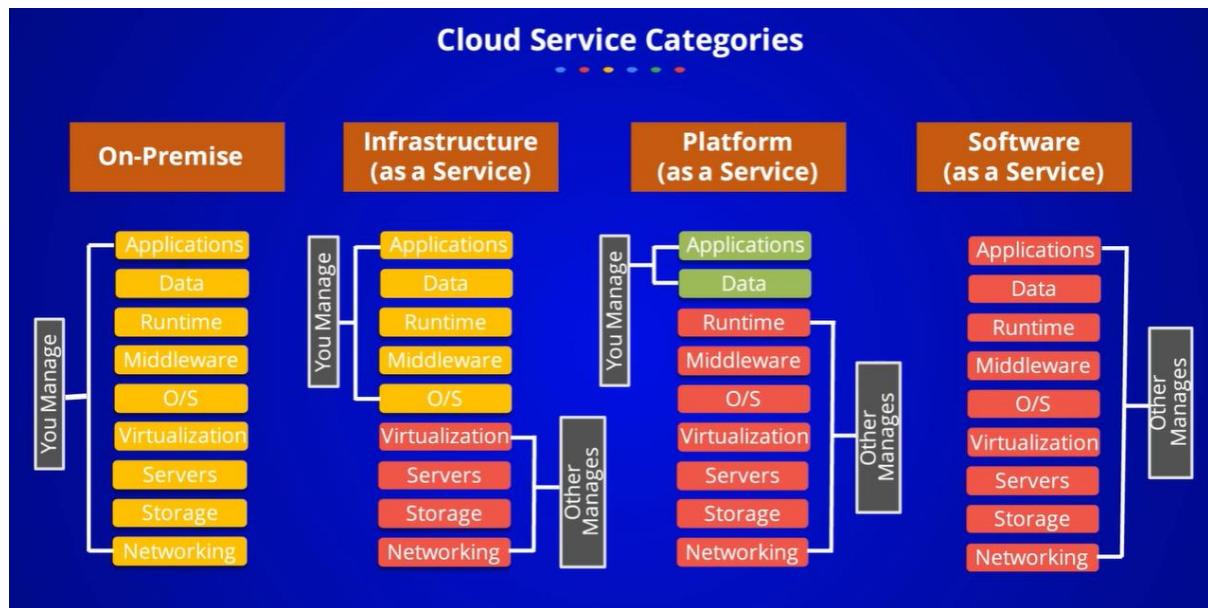


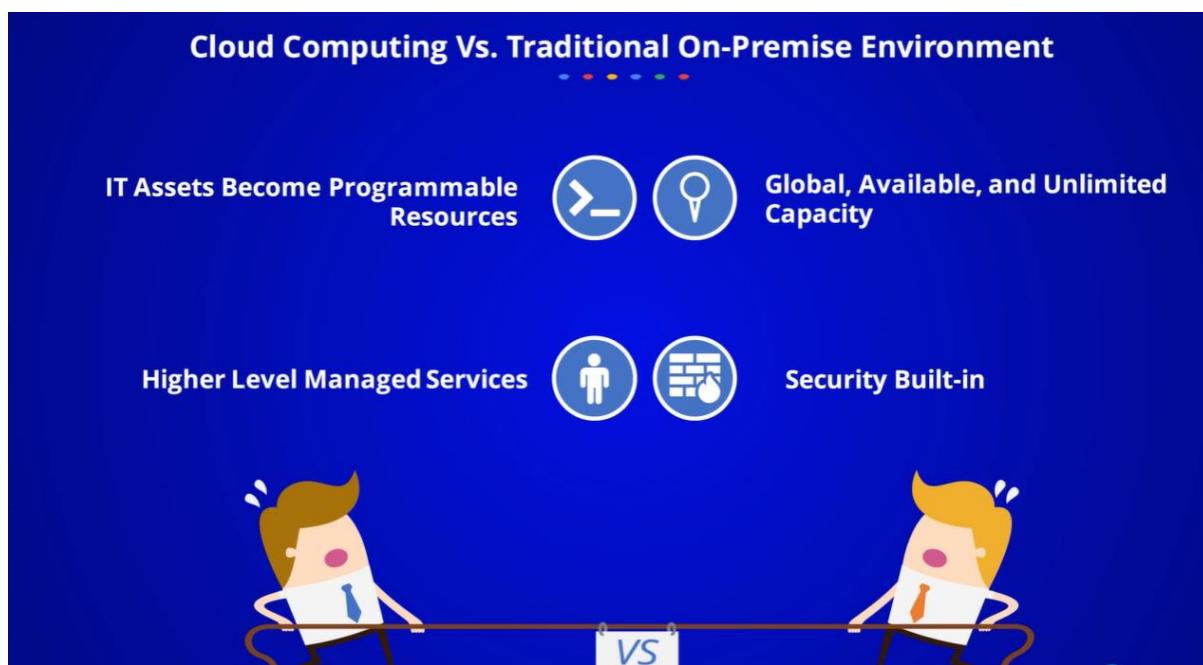
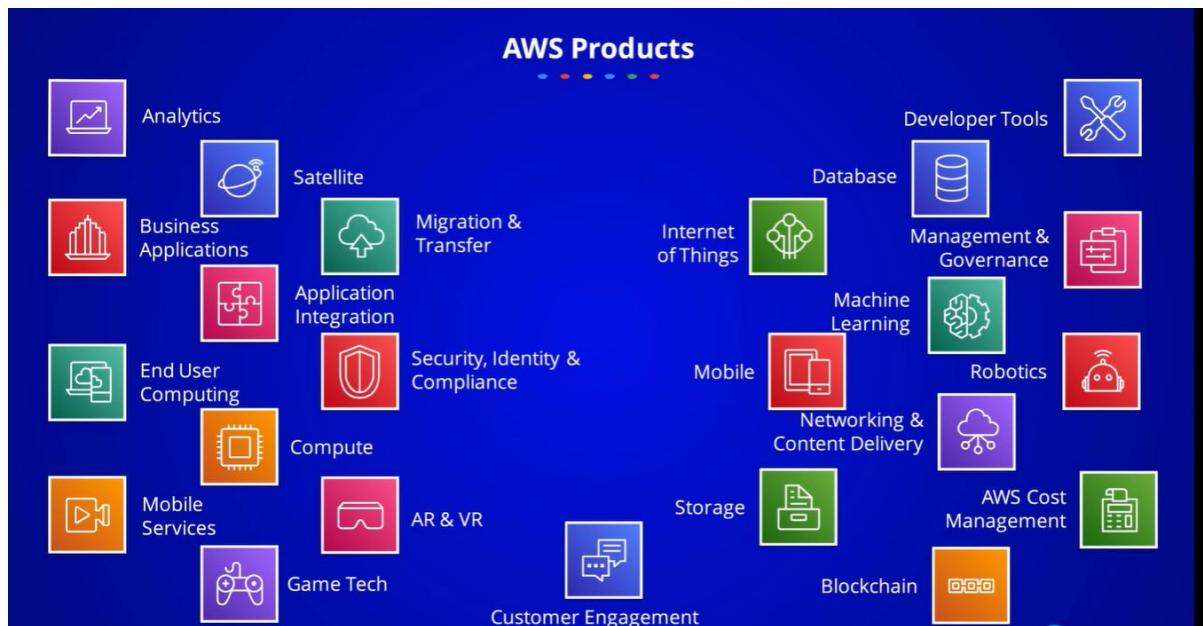
PaaS (Platform as a Service)

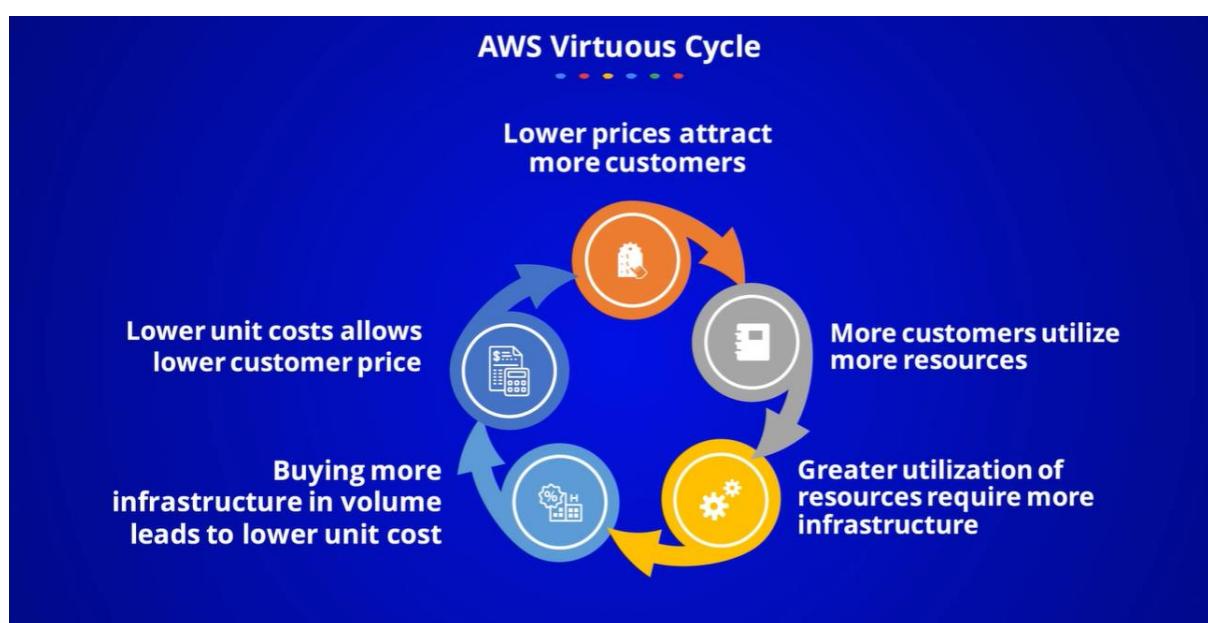
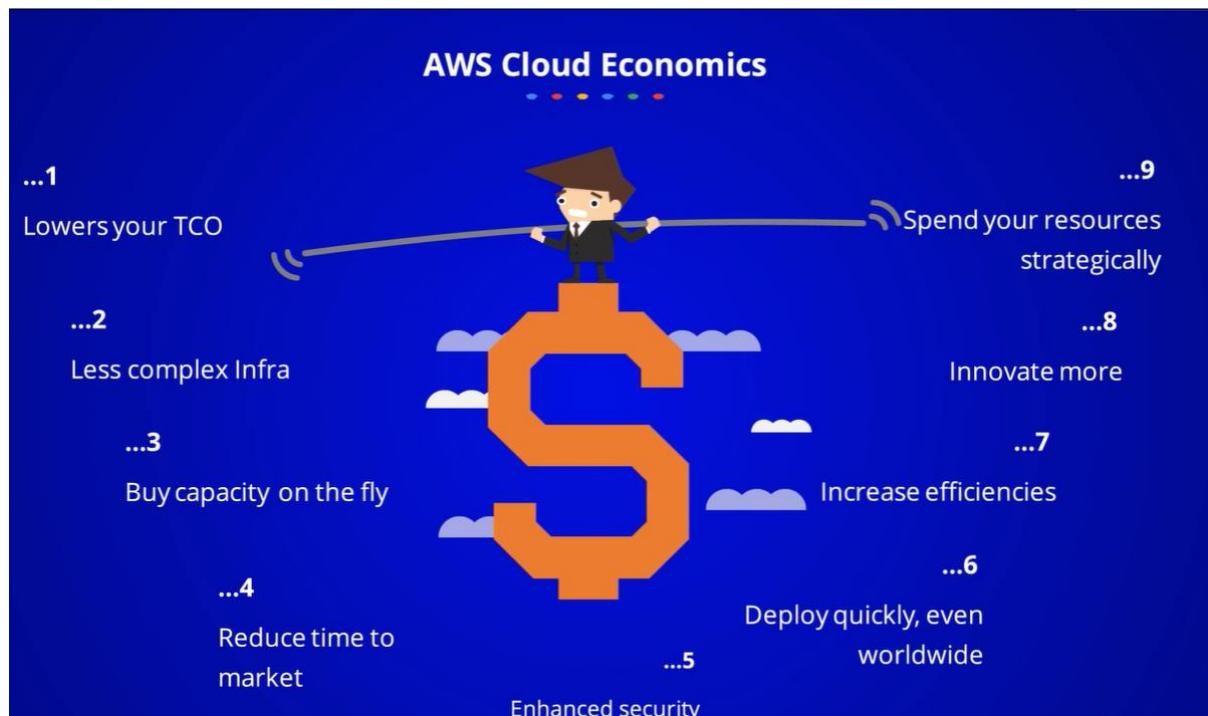


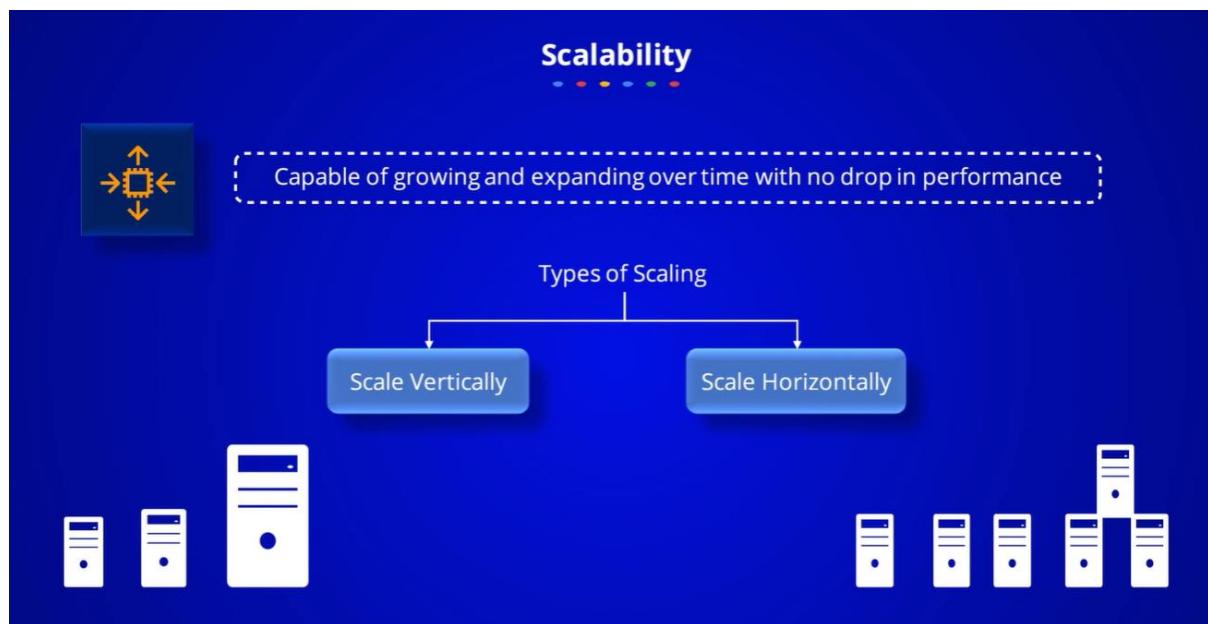
SaaS (Software as a Service)











Scale Vertically : adding more servers and added high ram and storage.

Scale horizontally: scale adding more machine to pool of resources to exiting machine.

Disposable Resources Instead of Fixed Servers

Auto Scaling permits EC2 instances to be disposed of when not utilized and new instances to be provisioned in seconds to satisfy the higher demand.

RDS database scales into several nodes and any node can be discarded without impact on DB availability



CloudWatch

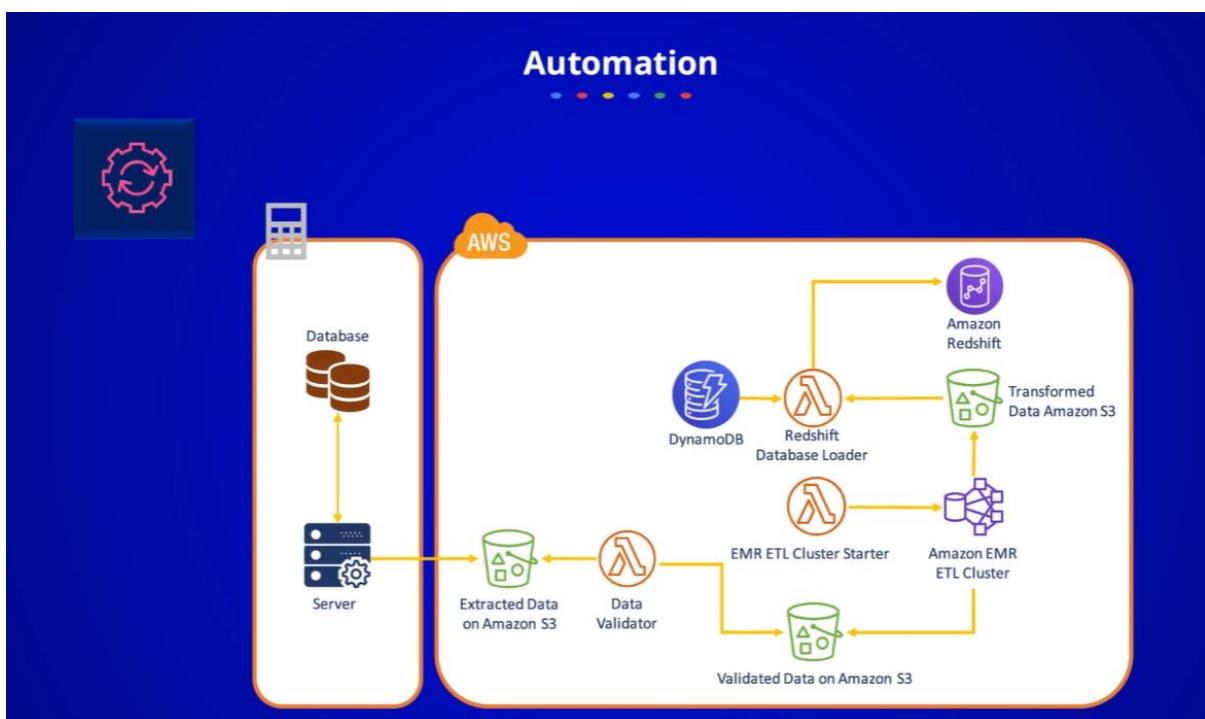


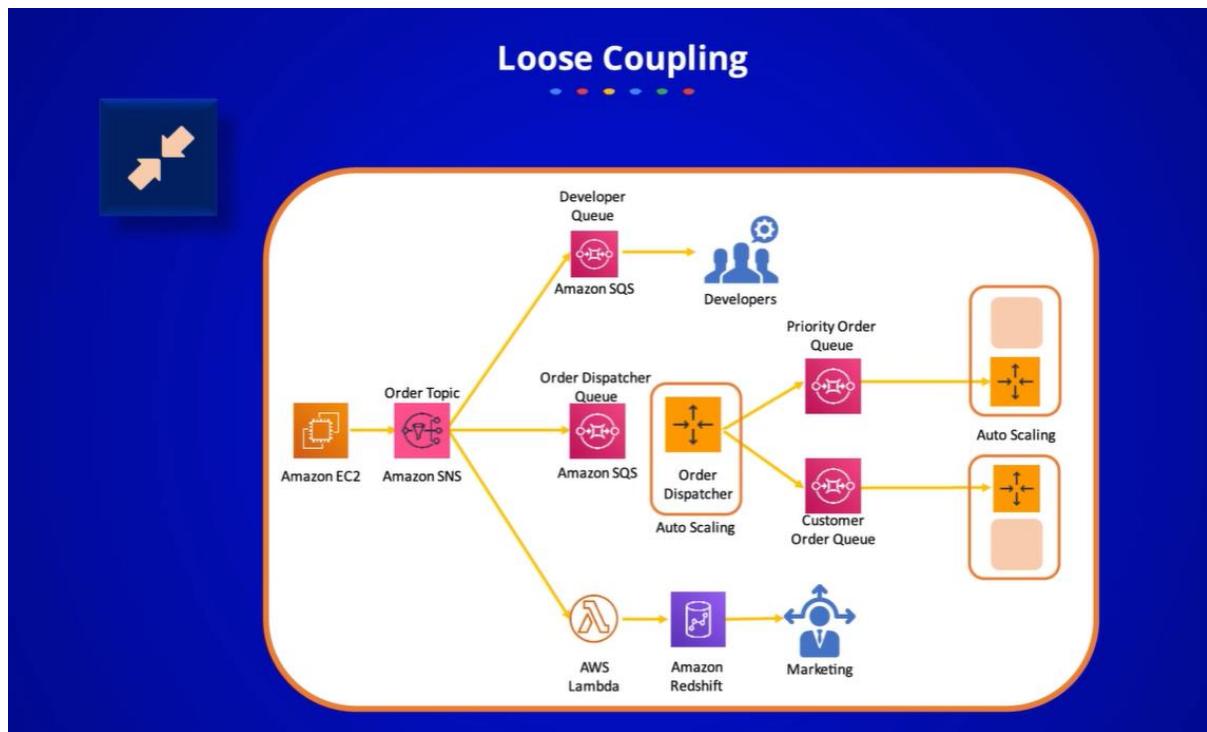
Instance crashes



Replacement automatically launches

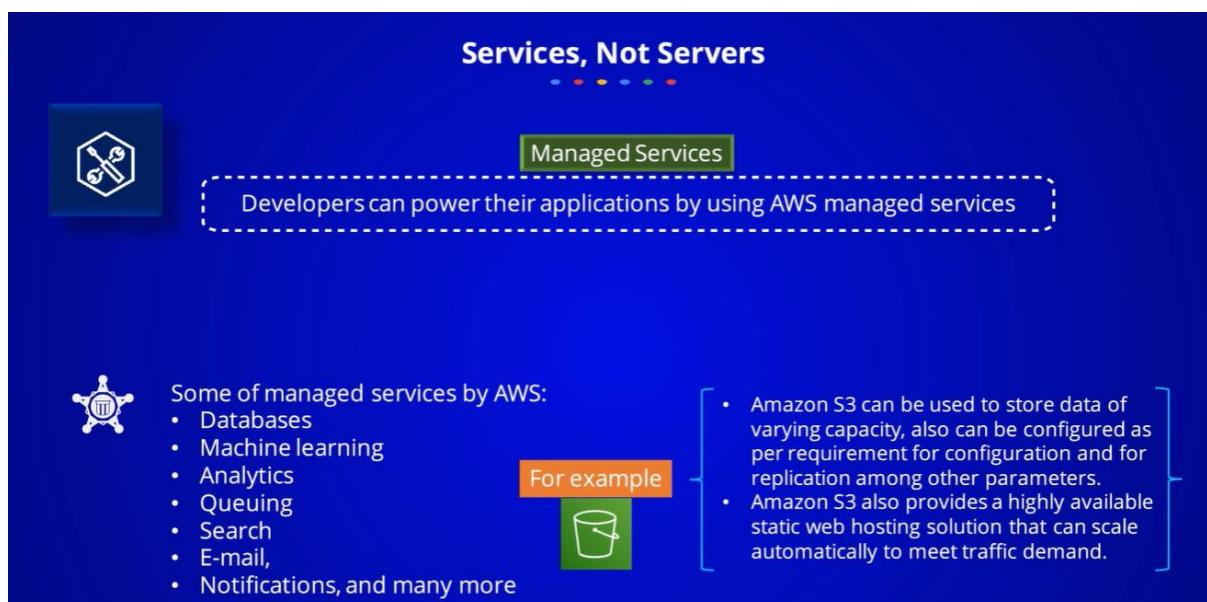
Automation





Loose coupling architecture.

Reduce dependences and for change reduce failer and not cascade to another compound .



Services, Not Servers

Server-less Architectures

Server-less architectures reduce the operational complexity of running applications.

Both Event-driven and synchronous services can be built without managing any server infrastructure.

For example

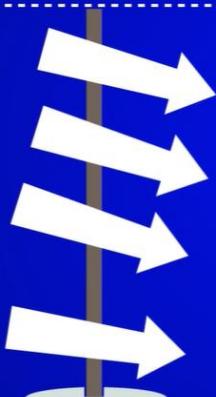


Your code can be uploaded to AWS Lambda compute service that runs the code on your behalf.

Removing Single Points of Failure



A system needs to be highly available to withstand any failure of the individual or multiple components (e.g., hard disks, servers, network links, etc.).



Introduce Redundancy

Automated Detection of Failure

Durable Data Storage

Automated Multi-Data Center Resilience

Caching



To store previously examined data, caching is used



Improves application performance



Increases the cost efficiency of implementation

Caching Types

Caching Types

Application Data Caching

Edge Caching

Application data can be stored in the cache for subsequent requests to:

- Improve latency for end users
- Reduce the load on back-end systems

Using Amazon CloudFront static and dynamic data can be cached at various edge locations across the world

Security



Businesses have legal obligations to keep client data secure

How to achieve security



Utilize AWS Features for Defense in Depth

- Use VPC for Isolation and WAF for Application protection



Shared Security Responsibility to AWS

- Security of the Cloud-AWS
- Security for the Data-Customer



Reduce Privileged Access

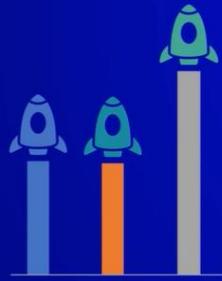
- To avoid a breach of security



Real-Time Auditing

- Continuously monitor and automate controls to minimize security risk exposures

Why AWS for Big Data-Challenges and Reasons



Big Data requires significant compute capacity that can vary in size.



Big Data on AWS meet your needs, without having to wait for additional hardware or being required to over invest to provision enough capacity.

Big Data on AWS lets you provision more capacity and compute in a matter of minutes, meaning that your big data applications grow and shrink as demand dictates

Why AWS for Big Data-Challenges and Reasons



Amazon S3



AWS Glue



AWS IoT



Lets you store data, orchestrate jobs and lets connected devices interact.

Amazon Kinesis Data Firehose



Lets you load streaming data continuously

AWS Database Migration



Lets you migrate data easily. AWS

AWS Direct Connect



Lets you have scalable private connections between On-premise and Cloud.

Databases in AWS

Benefits of AWS Database Services

 <p>Fully managed services AWS handles installs, patching, restarts</p>	 <p>Easy to scale Grow as you need</p>	 <p>Pay only for what you use No up-front cost</p>
		 Amazon S3  Amazon EC2  Amazon Data Pipeline  Amazon VPC  Amazon SNS  Amazon CloudWatch Designed for use with other AWS services

Relational Databases

Non-Relational Databases

Data Warehousing in AWS

Databases in AWS

AWS Database Services

Scalable High Performance Application Storage in the Cloud

Deployment & Administration

Application Services

Compute Storage **Database**

Networking

AWS Global Infrastructure


Amazon RDS


Amazon DynamoDB


Amazon Elastic Cache


Amazon Redshift


Amazon Database Migration Service

Databases in AWS



Amazon Relational Database Service (Amazon RDS) is used to setup, operate and scale relational database in cloud.



Automates time-consuming administrative task with cost efficiency and resizable capacity:

- hardware provisioning,
- database setup,
- patching and backups.



Allows to focus on following tasks for the application:

- **fast performance**,
- **high availability**,
- **security and compatibility**.

Databases in AWS

 Amazon RDS is available on multiple database instance types which one can choose from according to the specific application requirement- optimized for memory, performance or I/O.

There are 6 well known database engines available with Amazon RDS:



Amazon
Aurora



PostgreSQL



MySQL



MariaDB[®]

ORACLE



Microsoft[®]
SQL Server[®]

Databases in AWS

Relational Databases

Easy to administer

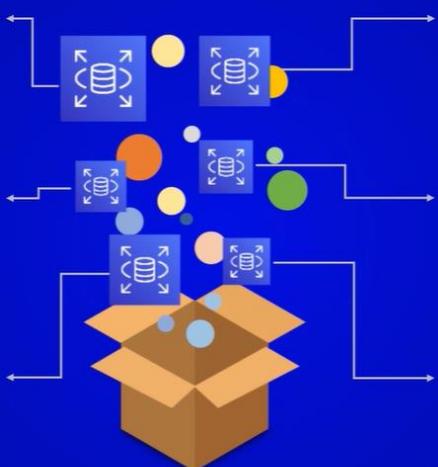
Fast

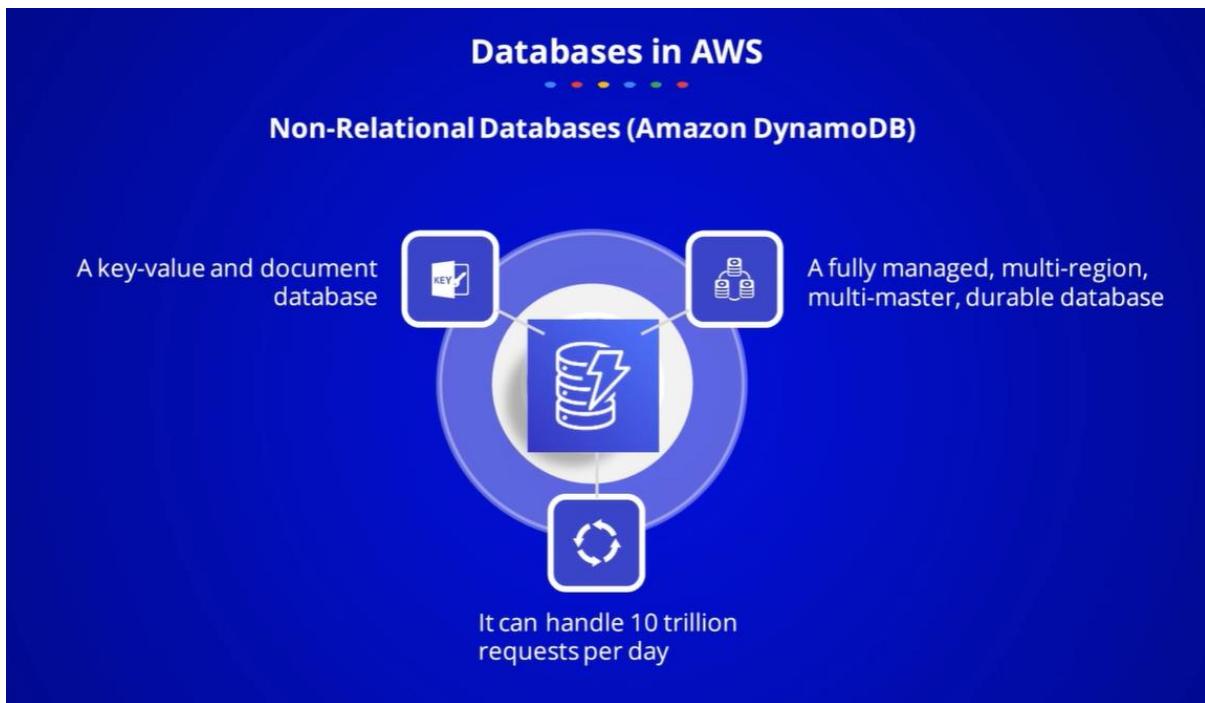
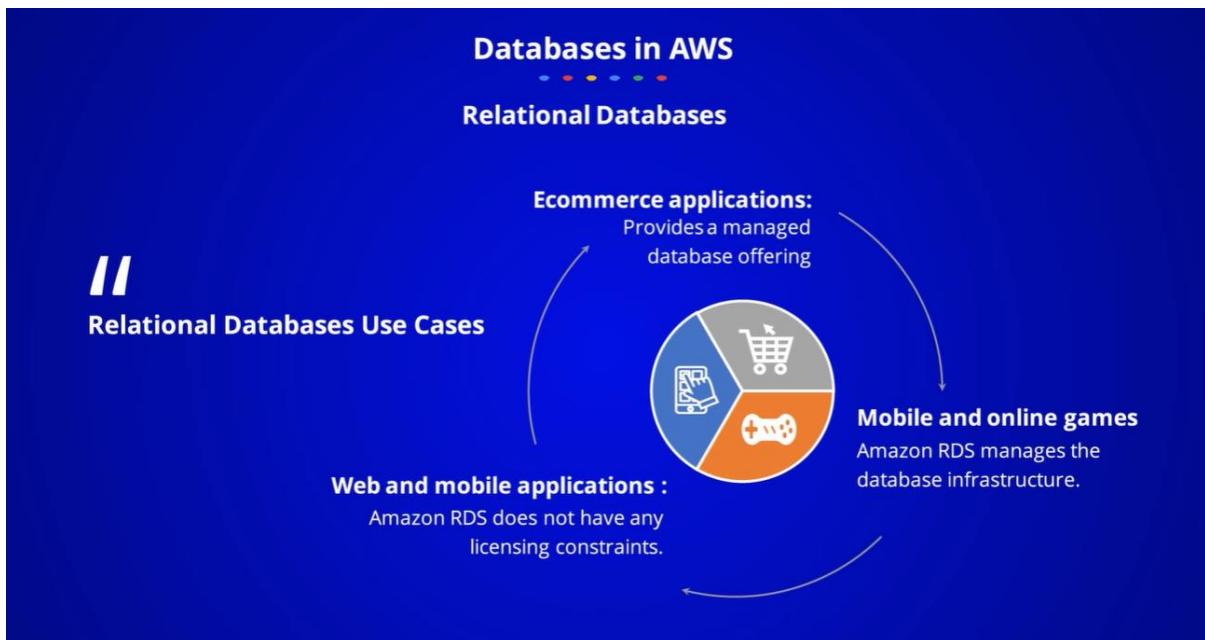
Highly scalable

Secure

Available and durable

Inexpensive





Databases in AWS

Non-Relational Databases (Amazon DynamoDB)

Benefits:



Performance at scale



You can build applications with virtually unlimited throughput and storage.



No servers to manage



DynamoDB is serverless with no servers to provision, patch, or manage and no software to install, maintain, or operate.

Databases in AWS

Non-Relational Databases (Amazon DynamoDB)

Use cases:

Ad-tech

A key-value store for storing various kinds of marketing data.

Gaming

Used to store the Game stats.

Retail

Uses common DynamoDB design patterns to deliver consistently low latency for mission critical use cases

Banking and Finance

Use fully managed services to increase agility, reduce time to market, and minimize operational overhead

Media and entertainment

Scales elastically to handle the high volume of load keeping latency low

Software and internet

able to handle internet-scale use cases and their requirements



Redshift is a relational database used for data warehousing

Data Warehousing in AWS

Amazon Redshift



Data warehouse:
Fast, fully managed and cost-effective



It gives you petabyte scale data warehousing and exabyte scale data lake analytics together in one service



It is up to ten times faster than traditional on-premises data warehouses.

Data Warehousing in AWS

Amazon Redshift

Key Features:



Maximize your ROI



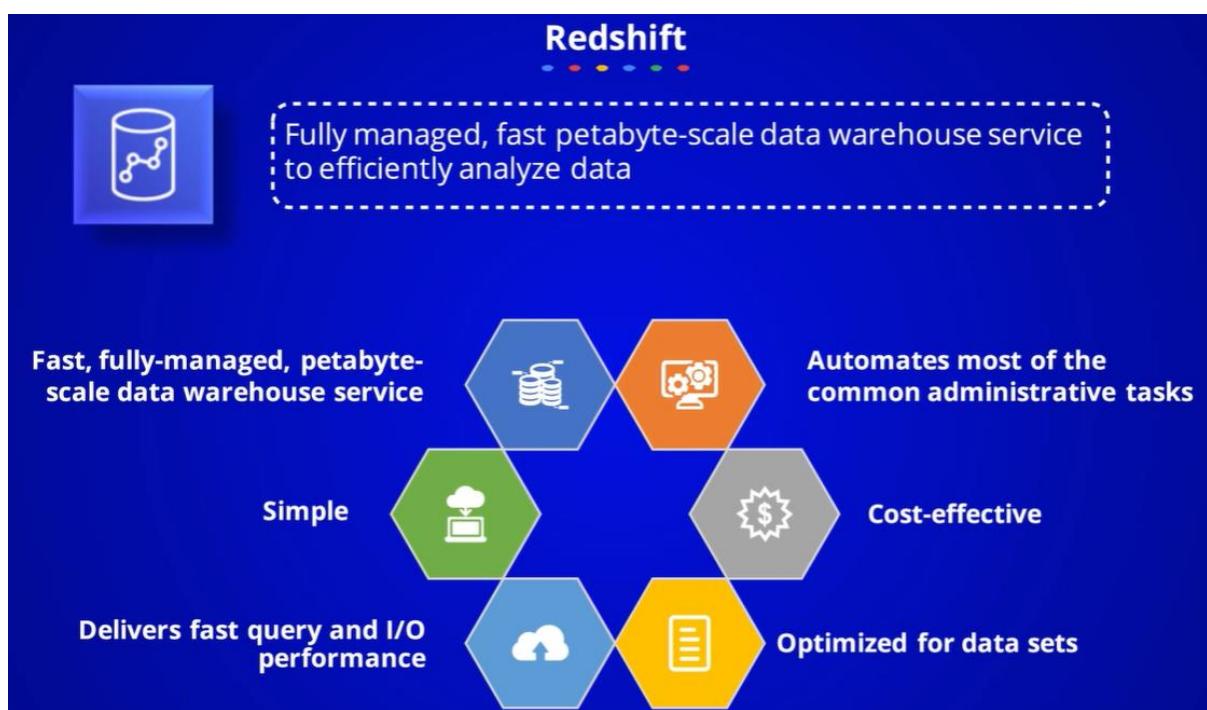
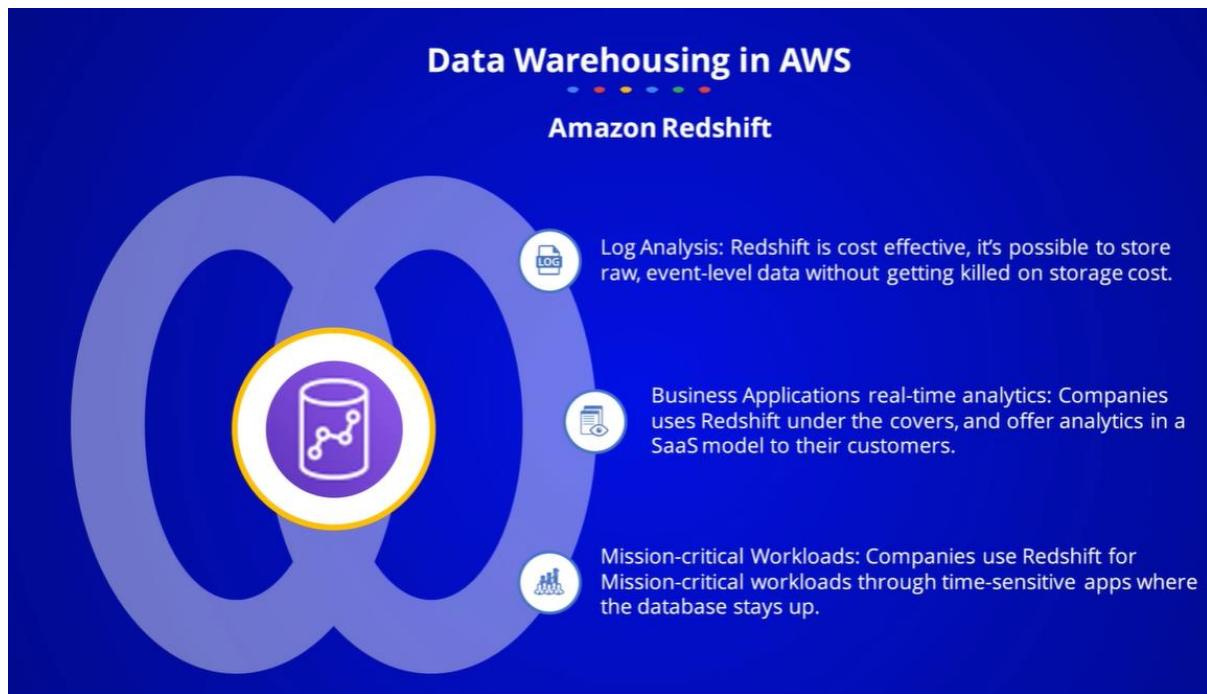
Amazon Redshift is a cost effective cloud data warehouse that costs less than 1/10th the cost of traditional data warehouses.



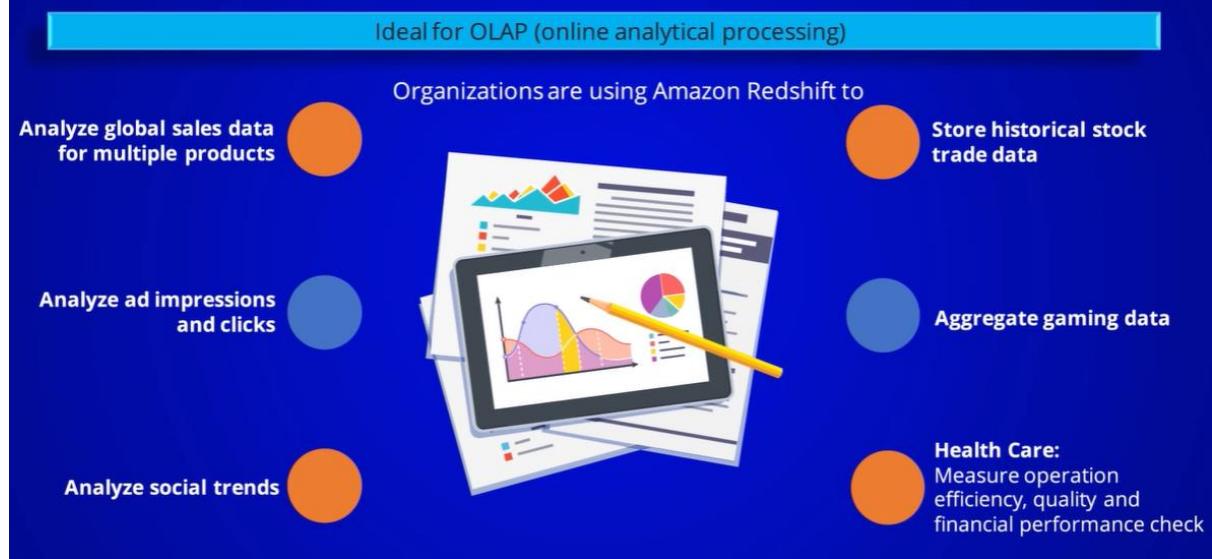
Setup, Deploy, and Manage with ease



Fast, fully managed, simple and cost-effective data warehouse

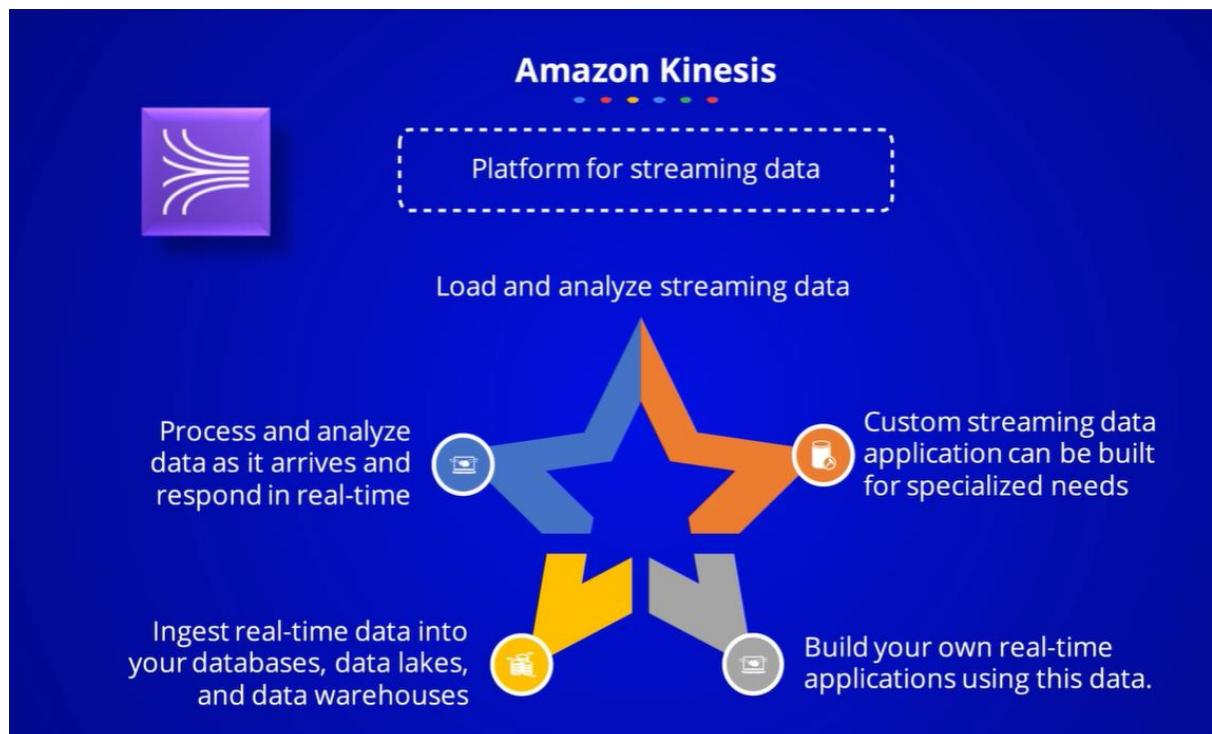


Redshift-Ideal Usage Patterns



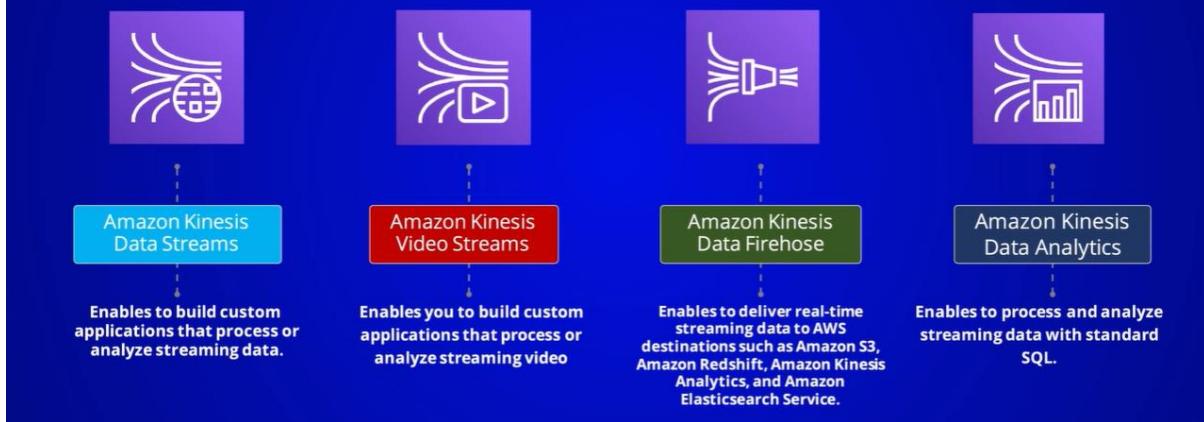
Kinesis

It's a platform for streaming data and to collect and process and analysis real time steaming data.

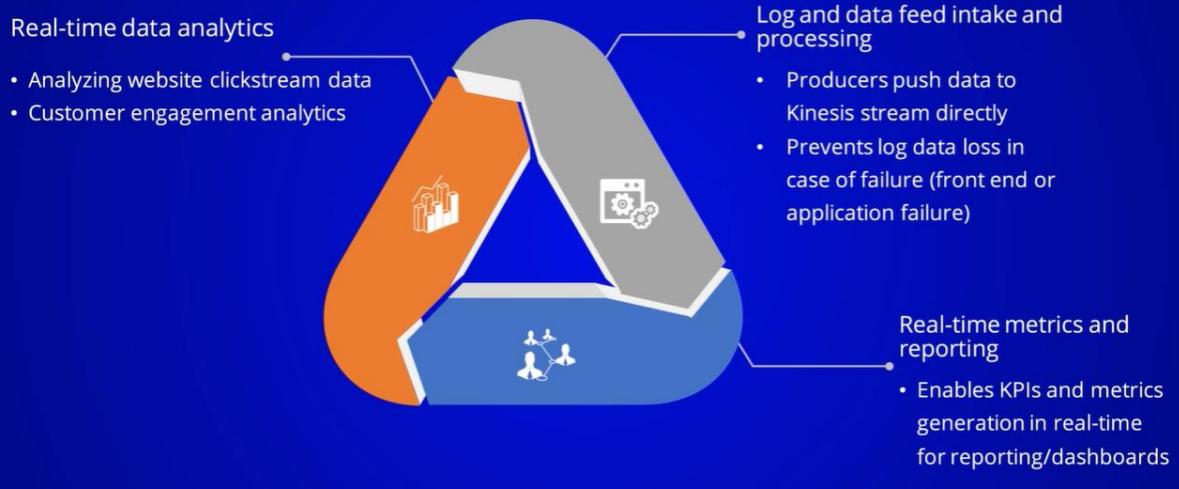


Amazon Kinesis

Currently there are 4 pieces of the Kinesis platform that can be utilized based on use case:



Amazon Kinesis-Ideal Usage Patterns



Amazon EMR



Highly distributed computing framework to process and store data quickly, easily and in cost effective manner.

Uses Apache Hadoop to distribute your data and processing across a resizable cluster of Amazon EC2 instances



Amazon EMR does all the work involved with provisioning, managing, and maintaining the infrastructure and software of a Hadoop cluster.

Amazon EMR-Ideal Usage Patterns



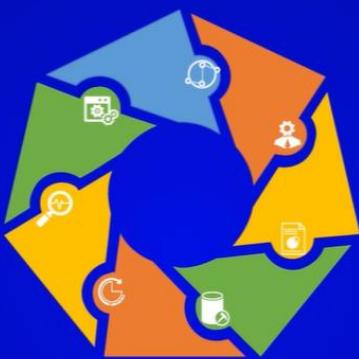
Log processing and analytics



Genomics



Ad targeting and click stream analytics



Large extract, transform, and load (ETL) data movement



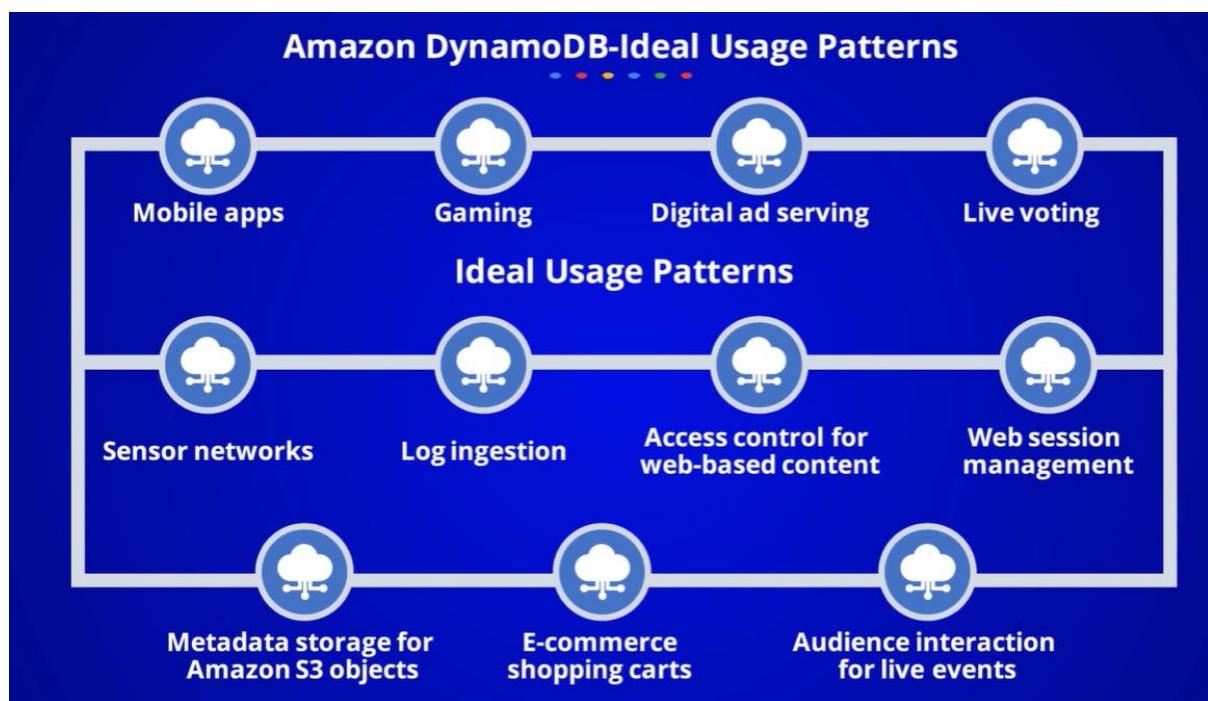
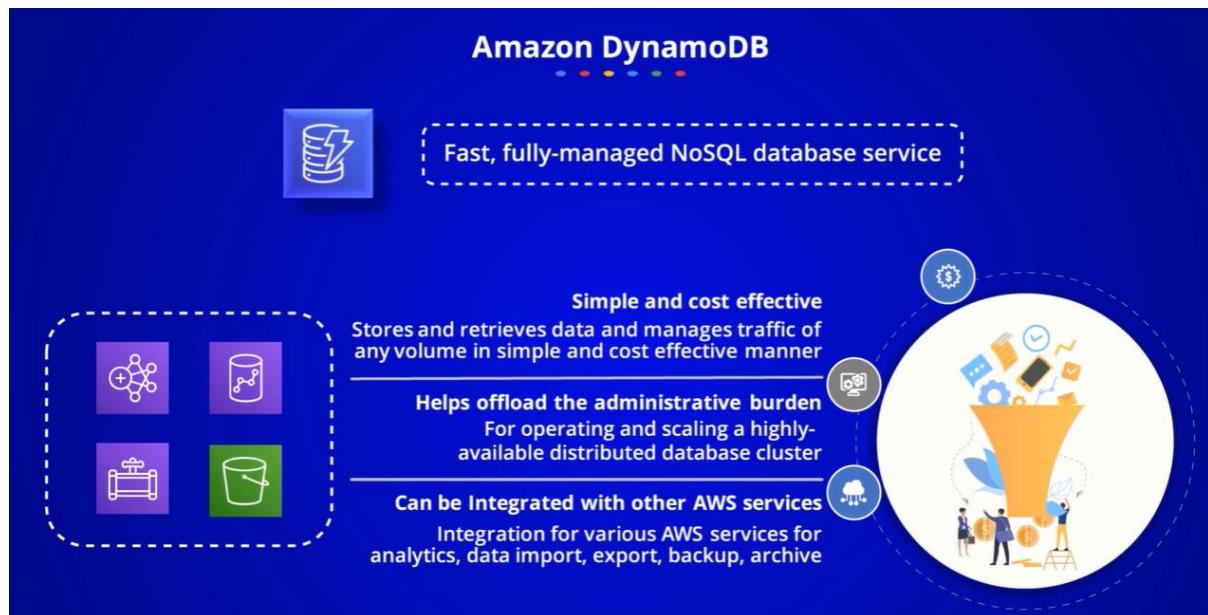
Risk modeling and threat analytics



Predictive analytics



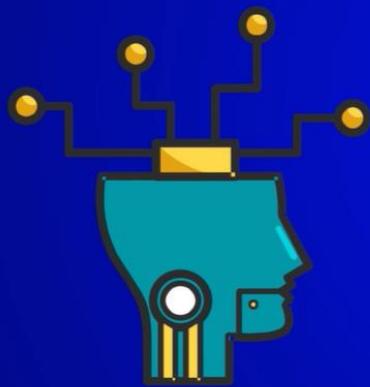
Ad hoc data mining and analytics



Amazon Machine Learning



Amazon ML is a service that makes it easy for anyone to use predictive analytics and machine-learning technology.



Provides visualization tool/wizards to create ML models
Learning of complex ML algorithms and technology is not required

Can create ML models based on data stored in AWS:
• S3
• Redshift
• RDS

Amazon Machine Learning-Ideal Usage Patterns

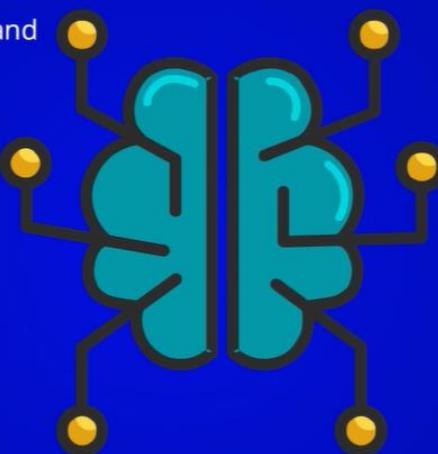
Forecast product demand

Predict user activity

Personalize application content

Listen to social media

Enable applications to flag suspicious transactions

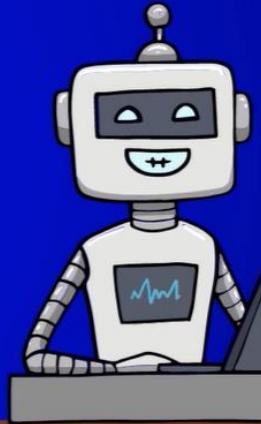


AWS Lambda

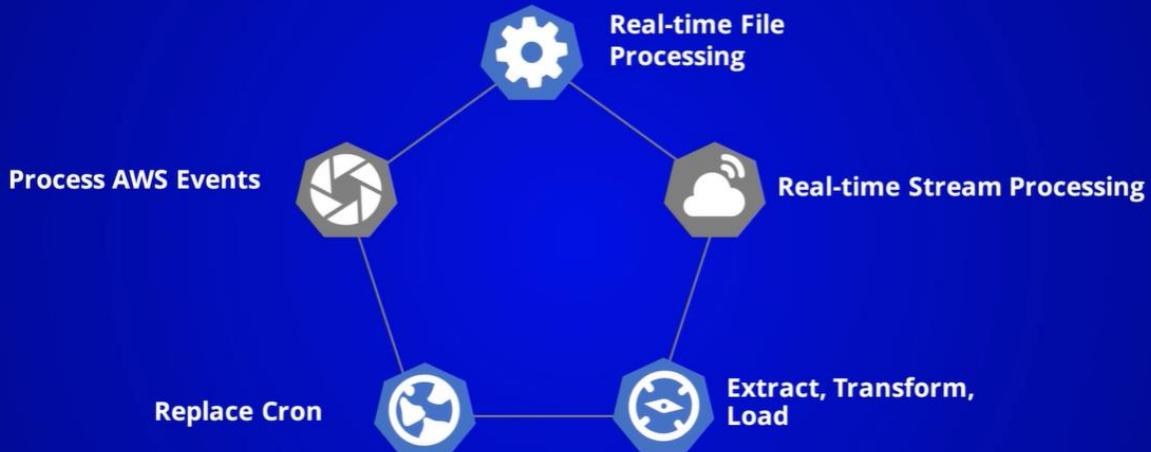


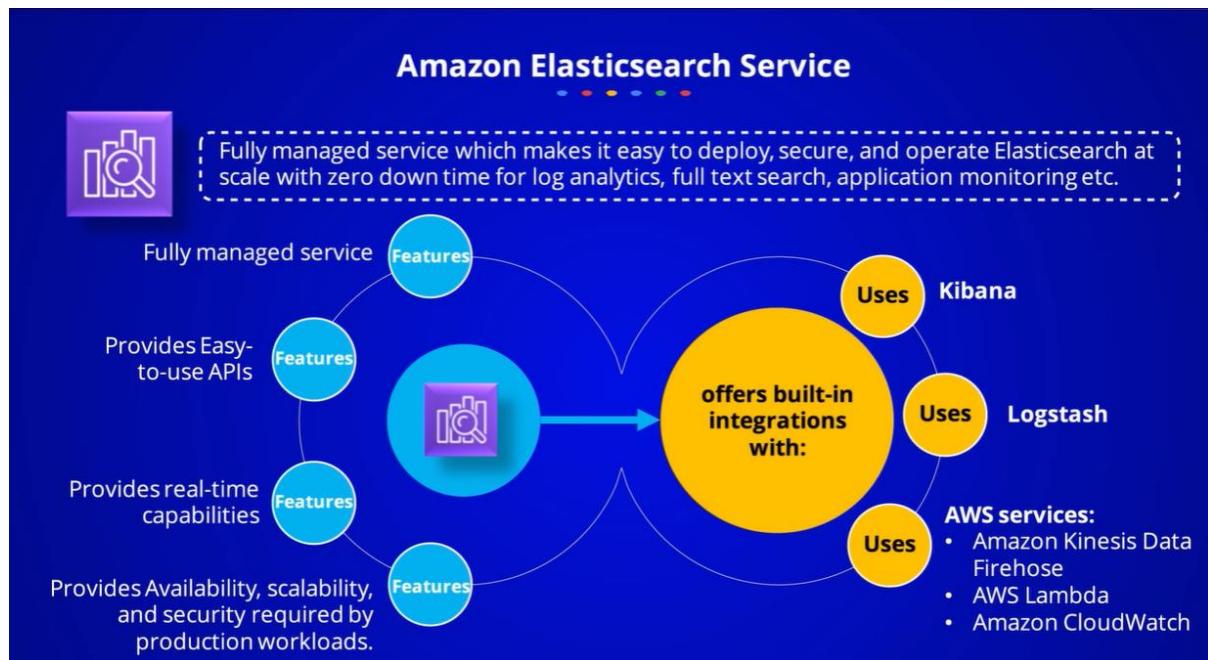
Let you run code directly on AWS without any provisioning/configuring

01. Pay only for the compute time you consume
02. Zero Administration: Run code on any application type or backend service with no administration
03. Just upload your code: Lambda takes care of everything else to run your code along with managing scalability and availability
04. Automatic triggering of code: Can be configured to trigger from other AWS services or from any other app (web/mobile)



AWS Lambda-Ideal Usage Pattern





Amazon EC2



Amazon EC2:

Amazon EC2, with instances acting as AWS virtual machines, provides an ideal platform for operating your own self-managed big data analytics applications on AWS infrastructure.



Almost any software you can install on Linux or Windows virtualized environments can be run on Amazon EC2 and you can use the pay-as-you-go pricing model.

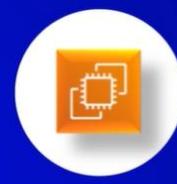


Amazon EC2

Ideal Usage Patterns



Specialized Environment



Compliance Requirements

Amazon EC2 provides the flexibility and scalability to meet your computing needs.

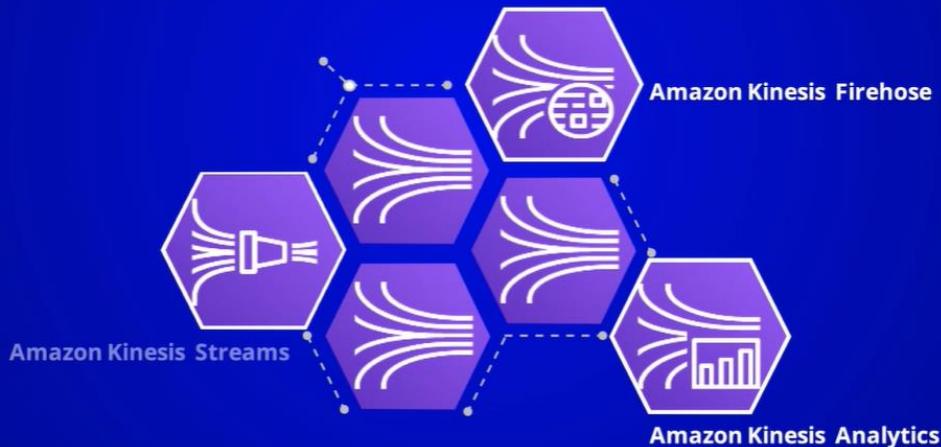
Certain compliance requirements may require you to run applications yourself on Amazon EC2 instead of using a managed service offering.

It's used to process collect and analyze real time and streaming data.

Kinesis Stream Introduction

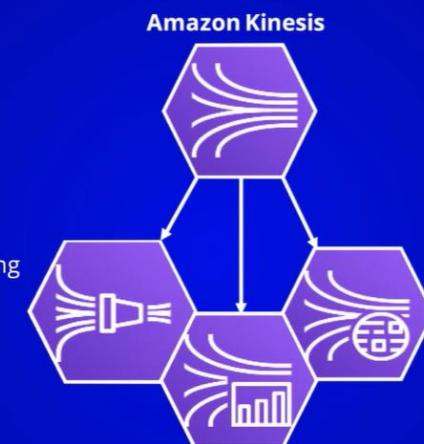
Amazon Kinesis is a service which is used to process, collect, and analyze real-time and streaming data.

Amazon Kinesis contains three different types of services:



Benefits of Amazon Kinesis Stream

Kinesis Stream is used for collecting and processing a large number of data records in real time.



Kinesis Firehose is used to load large streaming data into AWS services such as RedShift.

Kinesis Analytics is used to process and analyze a large amount of data in real time

Amazon Kinesis Stream

Amazon Kinesis is used to collect and process a large stream of data in real time.
The process records can then be sent to the dashboard.

Use Cases

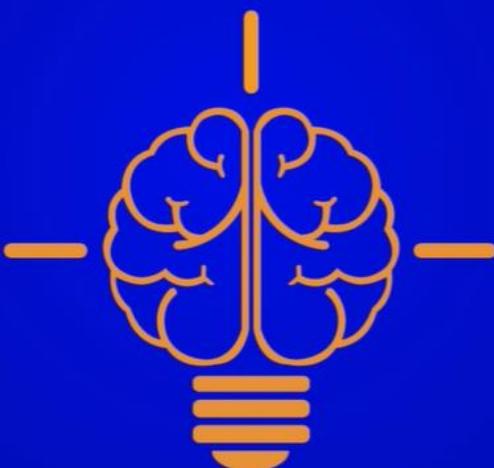


- Logging
- Real-Time Metric and Reporting
- Real-Time Data Analytics and Complex Streaming Processing

Benefits of Amazon Kinesis Stream

Fully managed

Real-time — Scalable



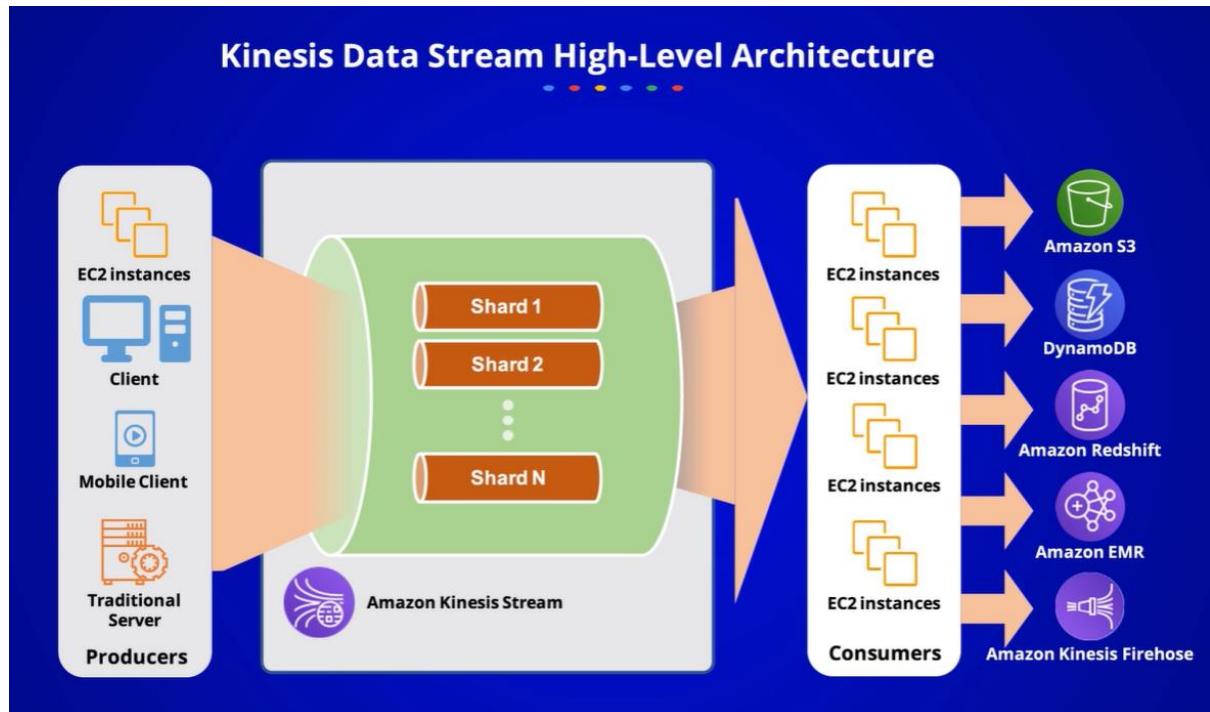
Logging

Real time metric and reporting

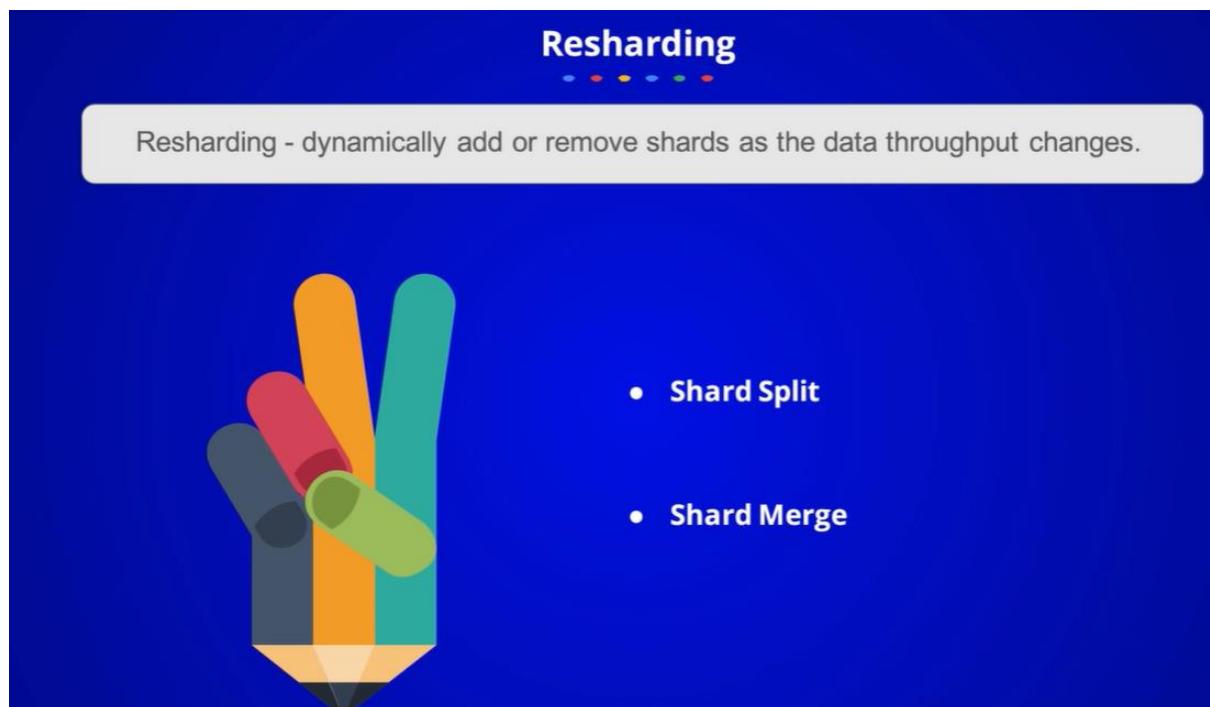
Real time data analytics and complex streaming processing.

Real time, full manage scalable.

- Producer App: add record into streaming.
- Consumer App: read data after the stream and record reading less then one sec.



Producer continuously push the code to kinesis stream, and consumers such as consumer app running consumers app to store in s3, dynamodb, redshift and emr, firehose.

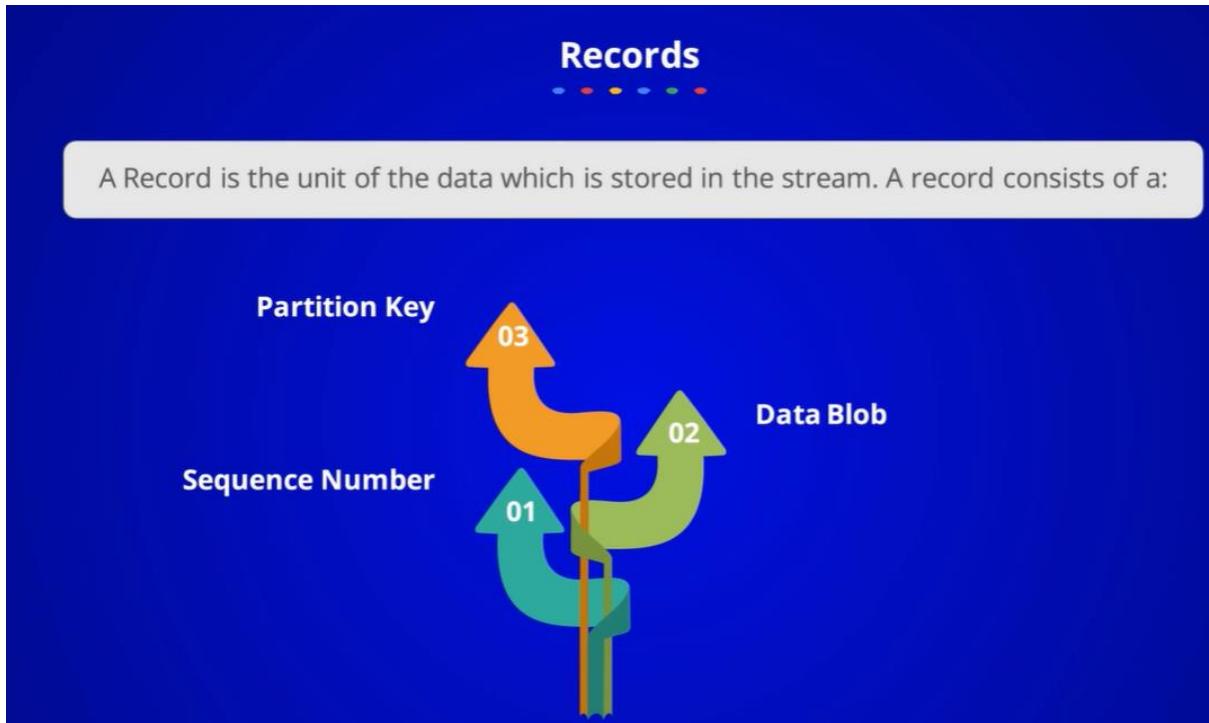


Two types of reshards

Shard split :

Shared split (single shared divided shared.)

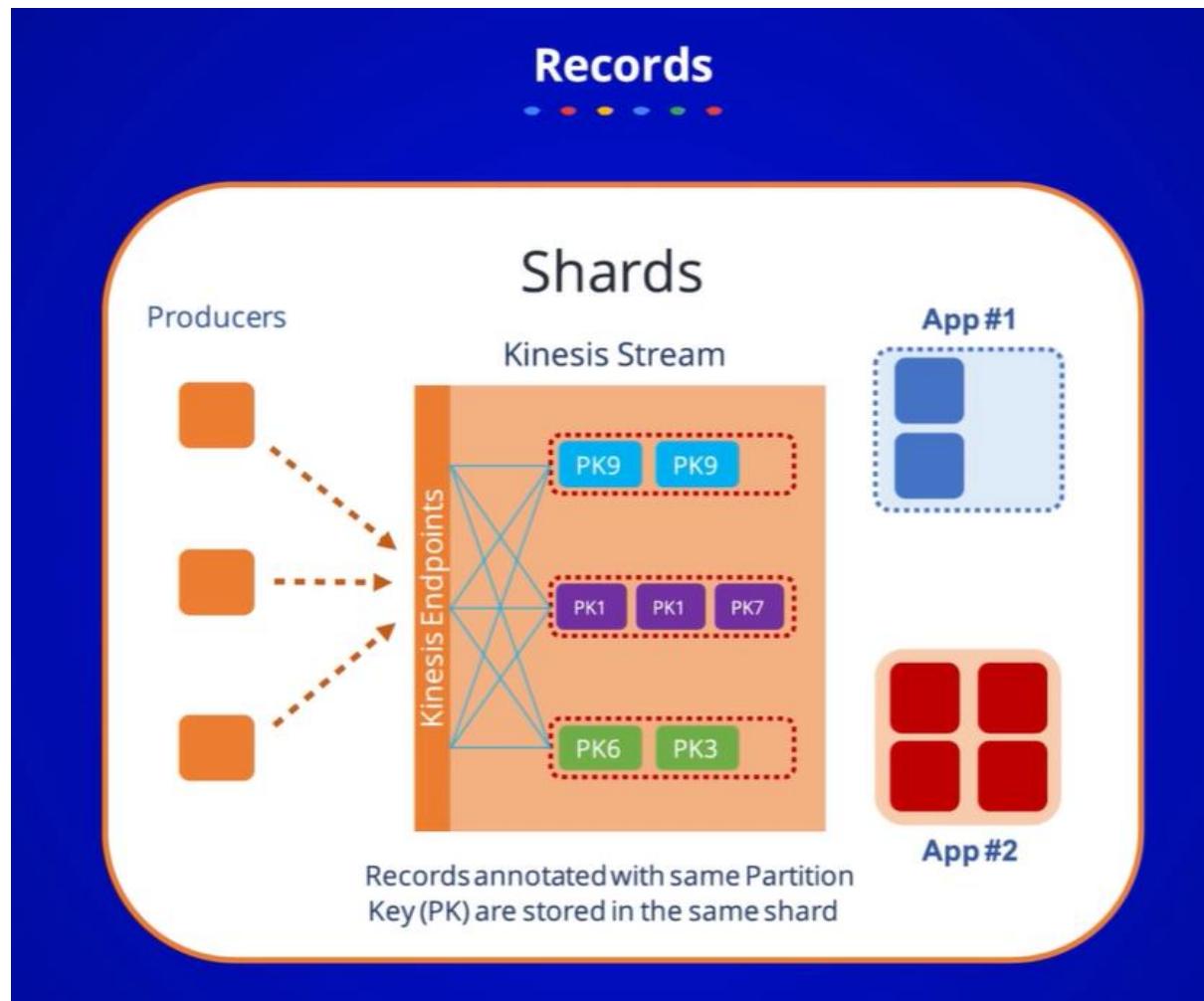
Shard merge(Shared merge two shared in single shard)



- Partition key(256 characters) is used stream distributed data across shards.

Kinesis data stream separated records belong to stream into multiple shards used partition key each associated data determine records which data record belong

- Sequence number:increase over time and specific shard within team.
- Data blob: size equal Or less then 1 Mb.



We have three shards contain difference partition keys in each shards and orgized are belong to sequence no is unique identifiers for records insert into shards.

Data blog: actual data data producer add to stream and max size 1Mb.

- Retention period: max duration after which data added into stream is expired
- Default retention period is 24 hrs by defaults and it can be increased to seven day if required.

Data Producer

Data Producer puts the data records to Kinesis Stream, and this could be something like the web server, a mobile client.



Amazon Kinesis Stream API

Amazon Kinesis Producer Library (KPL)

Amazon Kinesis Agent

Amazon Kinesis Stream API

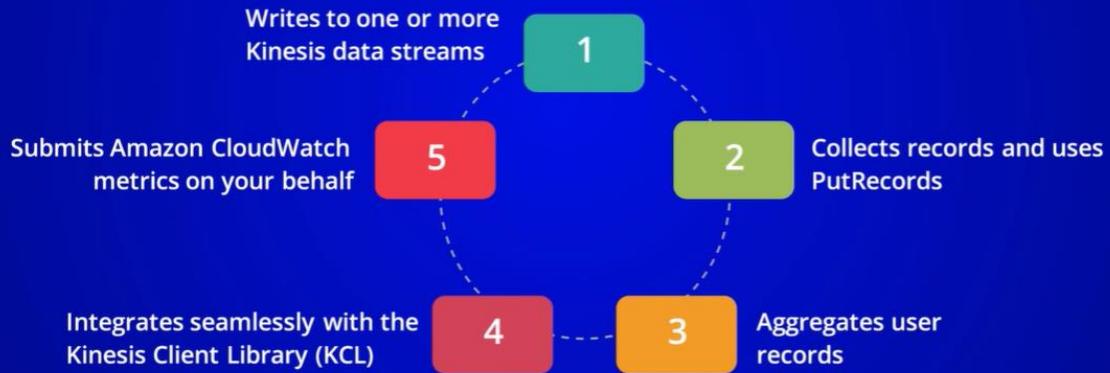


- It is a managed service that scales elastically for real-time processing of streaming big data.
- You can develop producers by using the AWS SDK for Java.
- PutRecord (Sends Single data record to the stream at a time.)
- PutRecords (Sends Multiple data records to the stream at a time.)

EXAM TIP: You should use PutRecords for most of the applications because it will give you a higher throughput for your data producer.

Amazon Kinesis Producer Library

The KPL is the configurable library which is used to create producer applications that allow the developers to achieve high-write throughput to a Kinesis Stream.



KPL Batching

Performing a single action on multiple items refers to as Batching. This reduces the effort of repeatedly acting on each item.

Batching helps you to use shard efficiently and also helps you with better throughputs.

KPL Batching

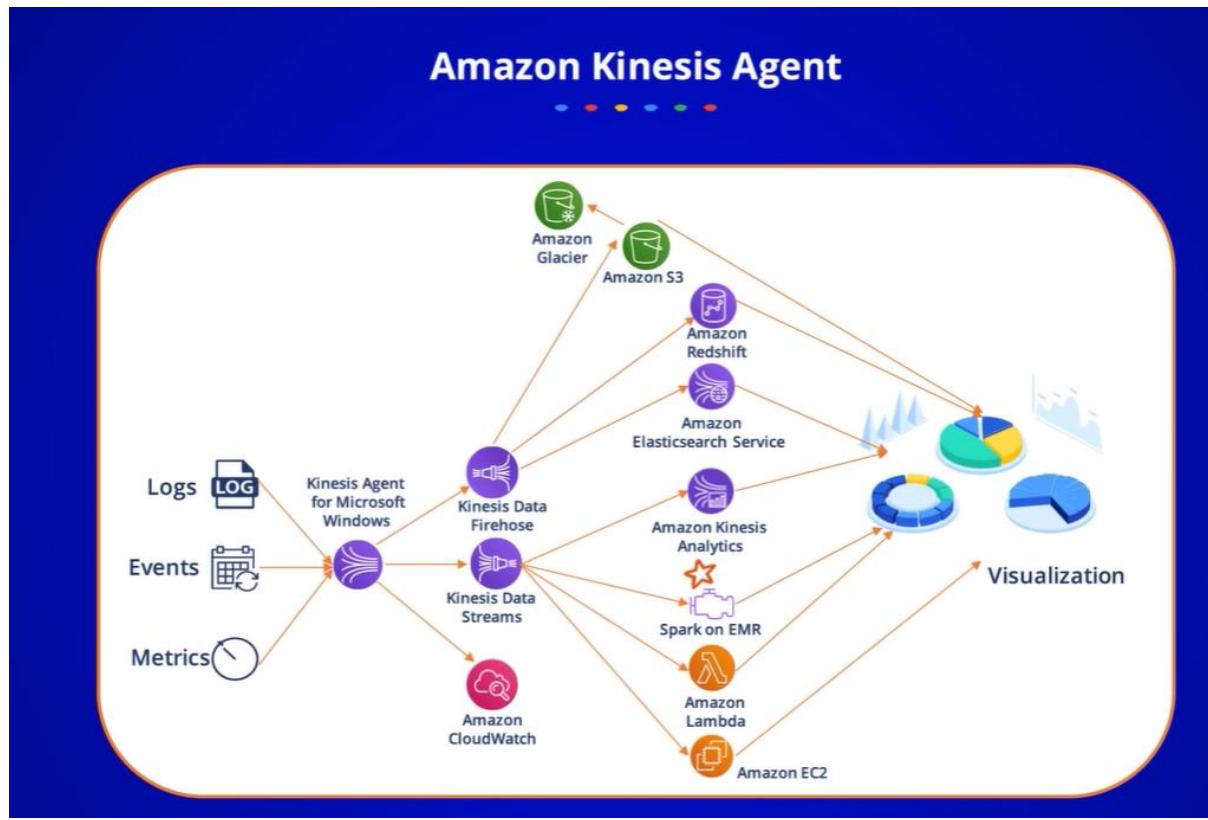
The KPL supports two types of batching:



- **Aggregation**
- **Collection**

Two types of batching.

- Aggregation (multiple record into single stream record)
- Collection(using API operation put the record send to multiple kinesis data stream record, one or more shards in kensis data stream.
By default data stream turn on.

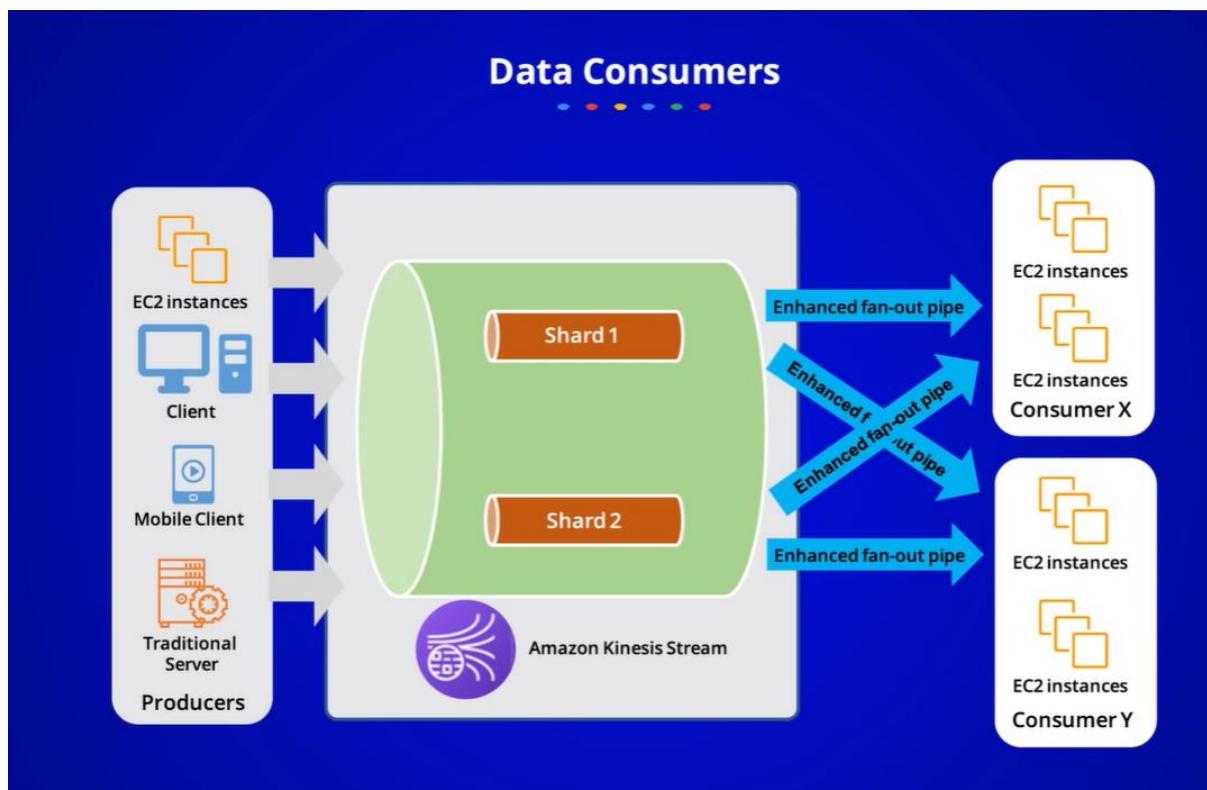
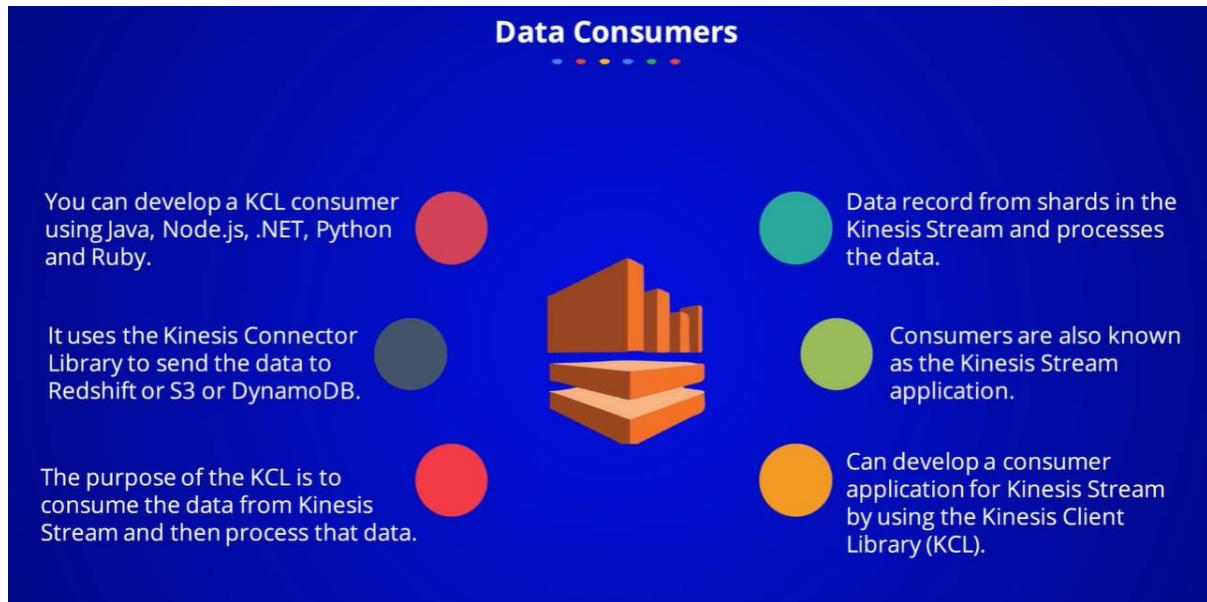


Glacier,s3, redshift and elastic search service.

Kinesis analytics, emr,lambda,ec2 finally go to visualization tools.

Data Consumers:

Consumer used to push data respectively from data stream.



Kinesis Stream Emitting Data to AWS Services



Consumers can emit data to S3, DynamoDB, Elasticsearch, Redshift, and EMR. Since multiple applications can consume data from the stream, you might have a case where some other data may be sent to services like RedShift, but the rest of the data may be archiving in S3 and use object lifecycle process to archive into the glacier.

Kinesis Stream Emitting Data to AWS Services

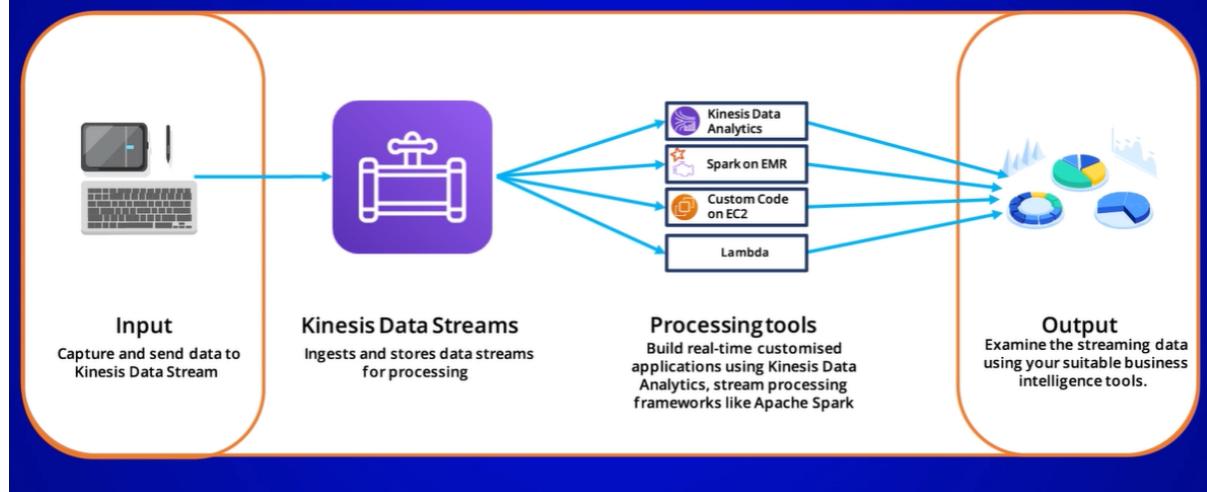
DynamoDB-Data in the stream could be sent to DynamoDB.

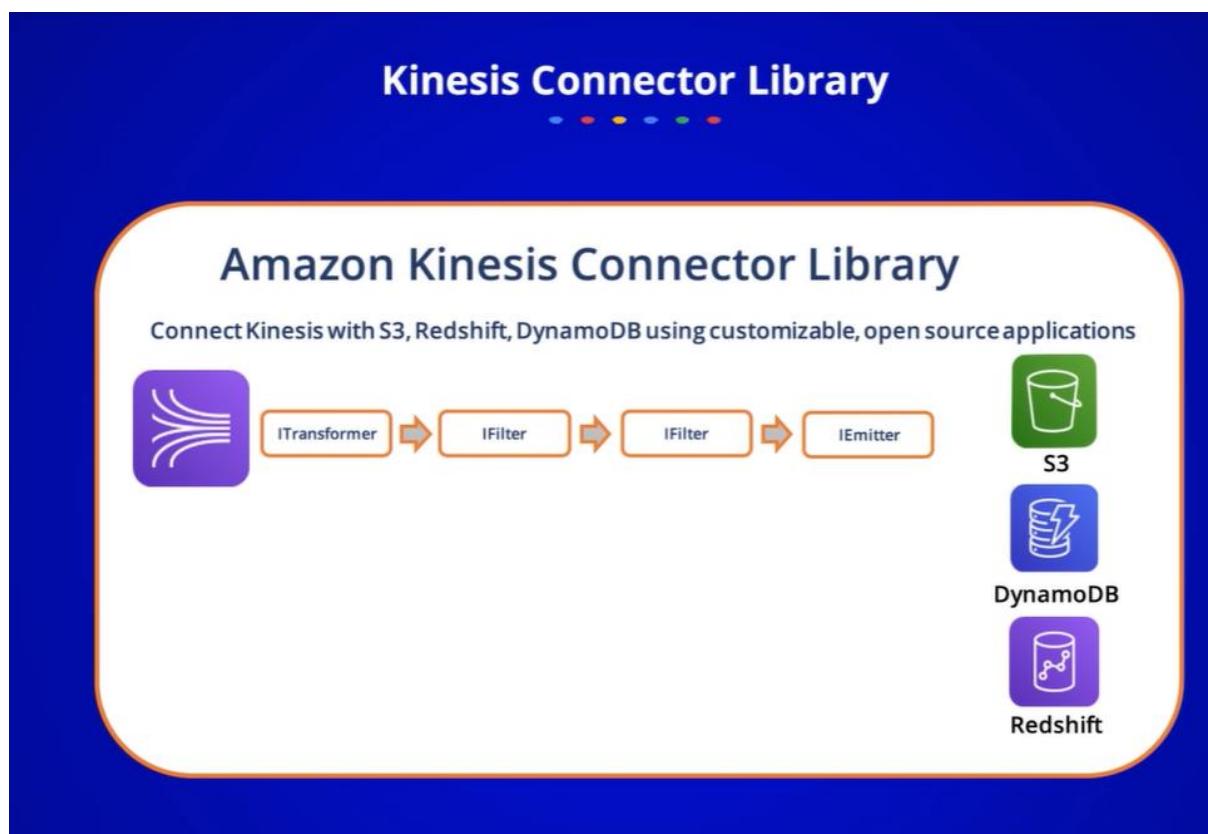
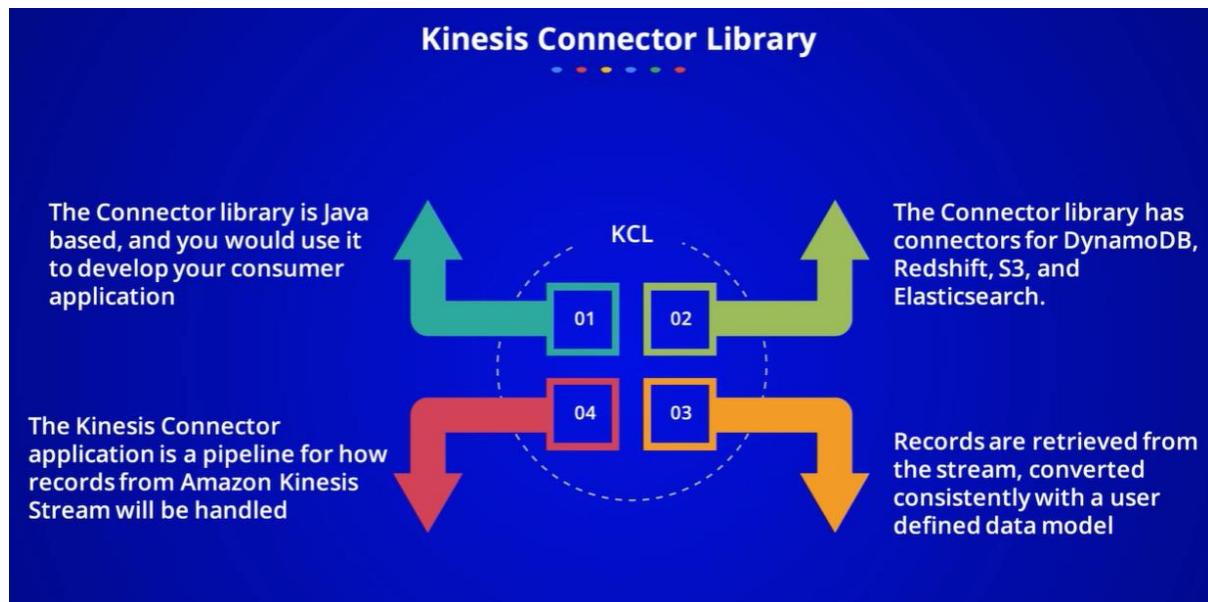
Redshift-If Data needs to be added to data warehouses send it to Redshift.



Elastic search-if you need to make the data from the stream searchable, then send this data to Elasticsearch using the KCL.

Kinesis Stream Emitting Data to AWS Services





Consumer can emit data to s3, dynamo db, elastic search, redshift and emr.

Since multiple app can consumer data from stream you might have case where some other data may be sent to services like redshift but rest of data may be archived in s3. And use object life cycle process to archive into glacier.

Use case

- Dynamo Db data in stream could be sent dynamo db

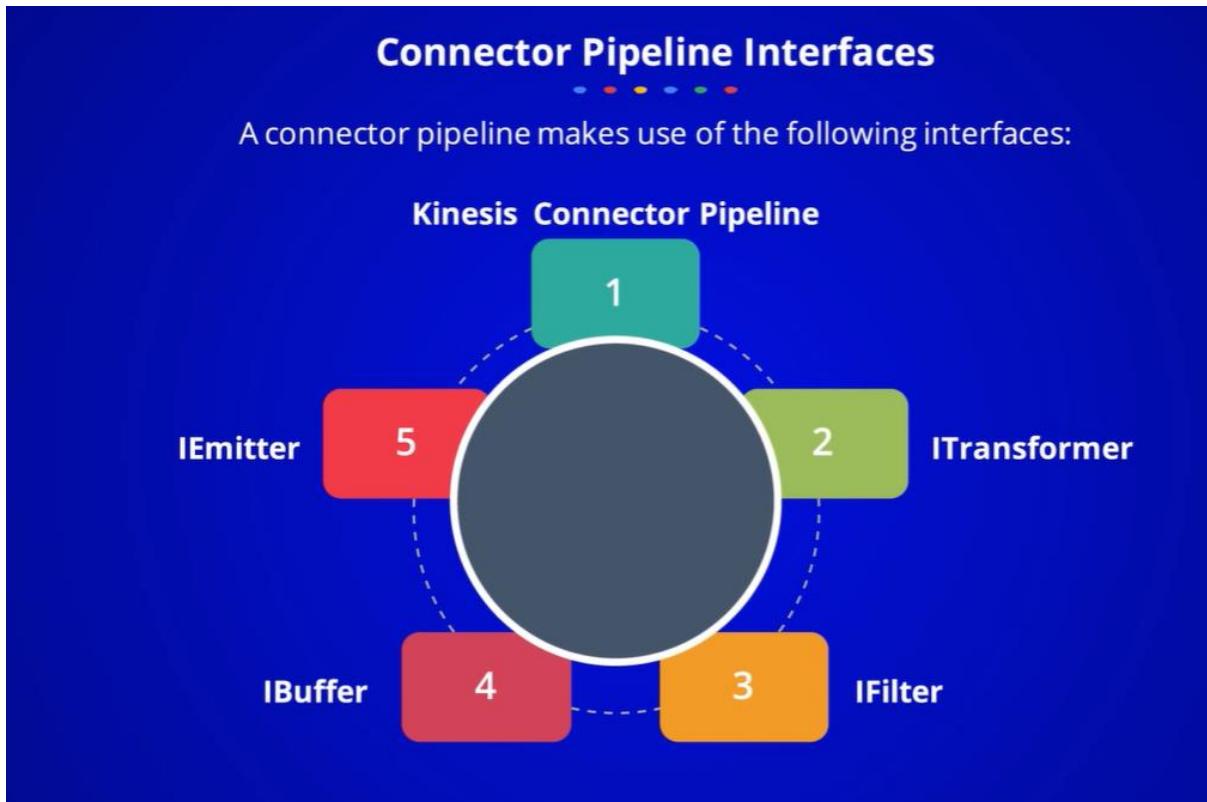
- Like gaming app and dashboard-> read data

Redshift : if data needs to be added to data warehouse send it to redshift.

Elastic search: if you need to make data from steam searchable then send this data to elastic search using kcl.

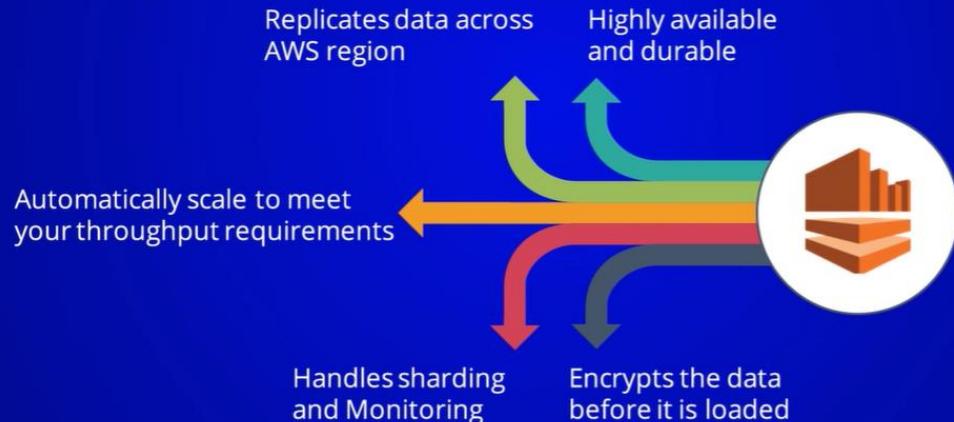
Capture and send data to kiness data stream

Kinesis data stream(ingests and stores data streams for processing)

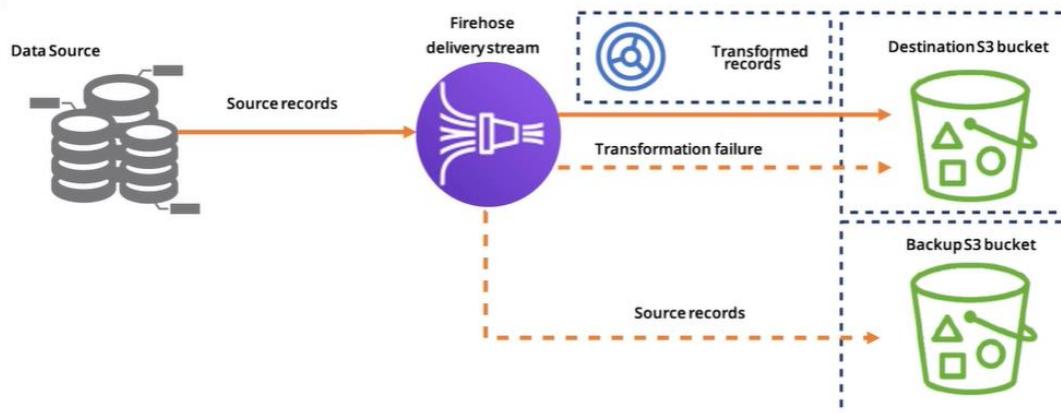


Kinesis Firehose

Amazon Kinesis Firehose is a fully managed service used to collect and load streaming data in near real-time into S3, Redshift, and Elasticsearch.



Kinesis Firehose

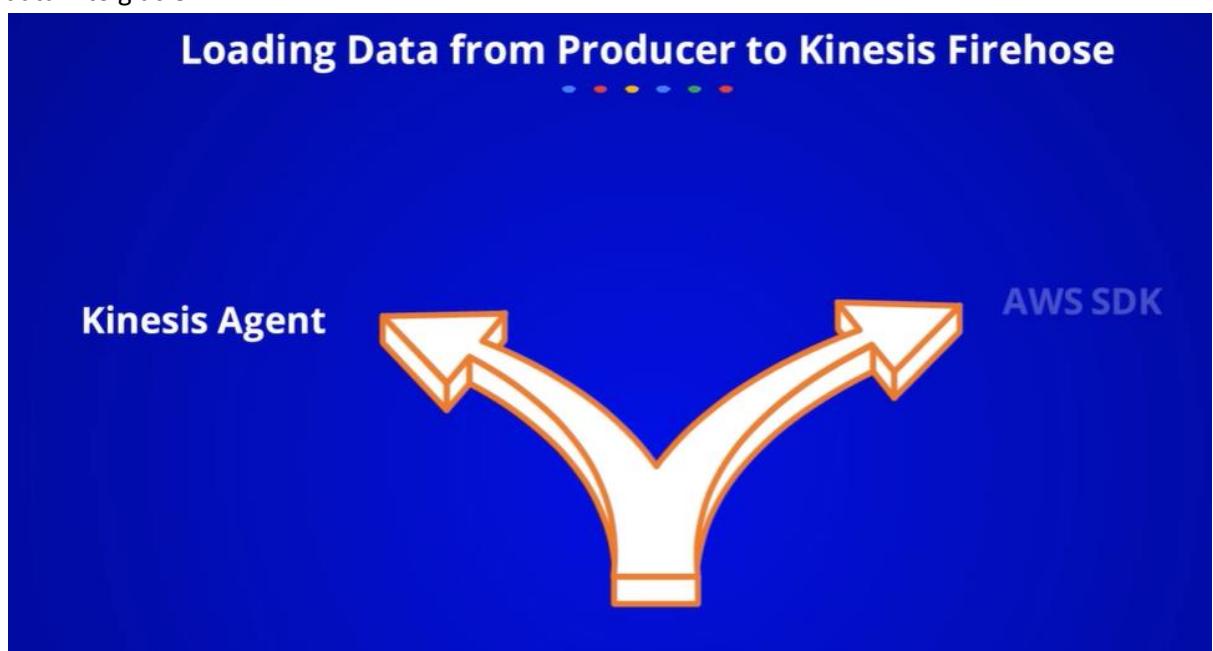


- No need application manage by resource
- We can config data source send source records to firehose delivery stream.
- It automatically delivery data to destination also configure transformed records.
- If you destination configure s3 delivery s3 bucket.



Web-server(producers): generate logs which send data to kinesis fire hose delivery stream in form of records.

- Each record can be as large as 10000kb
- For each delivery stream you will set buffer size or buffer interval.
- Firehose buffer incoming streaming data to certain size of certain period before delivering it to s3 or elastic search.
- Buffer size can be selected from 1mb to 128mb or buffer interval of 60 to 900 sec.
- Condition which is satisfied first it triggers data delivery to s3 or elastic search.
- From firehose delivery stream the data is then put into s3 and the using copy command.
- After data has been pulled into redshift, you can use the object lifecycle path to archive the data into glacier.



Two types
Kinesis agent
Aws sdk

Agents

- 1) Standalone java software application offers easy to collect data into data stream and stream kinesis firehose agent continuously set of file and seniers data into stream.
- 2) Allows you to collect and send data to firehose.
- 3) Can install it on linux based web server log server or data base web servers.
- 4) Will also help in cloudwatch metrics.
- 5) Can preprocess records from monitored files this involved multiple line in single line converting delimiters to json format and converting records to log format to json format.
File monitoring destination stream for data
- 6) Agent can ve monitoring multiple tree and write multiple streams it also preprocess covert records for send to stream or delivery stream.



Firehost agent have Two types of records

* **Put records.**

Single data record into amazon kiness firehose delivery stream.

***put record batch**

Use multiple record into delivery stream in single call it can achieve high through put then writing single records

Kinesis Firehose



- Don't write any code just create delivery stream and configure destination data
- Client data into stream using aws API calls and put records batch and data automatically sent to destination.
- When configure say stream Amazon s3 and firehose and data directly to s3.
- Amazon red shift (destination data shift to Amazon s3 and then red shift copy command execute load data to red shift).

Workout



Put and get records from kinesis data stream

Go to cli

Aws

Step1

```
aws. error. too few arguments
CHMB119868:Downloads samuelasirvat_r$ aws kinesis create-stream --stream-name Foo --shard-count 1
```

Step2

```
CHMB119868:Downloads samuelasirvat_r$ aws kinesis describe-stream --stream-name Foo
{
    "StreamDescription": {
        "KeyId": null,
        "EncryptionType": "NONE",
        "StreamStatus": "CREATING",
        "StreamName": "Foo",
        "Shards": [],
        "StreamARN": "arn:aws:kinesis:us-east-1:559483139917:stream/Foo",
        "EnhancedMonitoring": [
            {
                "ShardLevelMetrics": []
            }
        ],
        "StreamCreationTimestamp": 1572504543.0,
        "RetentionPeriodHours": 24
    }
}
```

Step3 checking stream active using previous cmd

```
CHMB119868:Downloads samuelasirvat_r$ aws kinesis describe-stream --stream-name Foo
{
    "StreamDescription": {
        "KeyId": null,
        "EncryptionType": "NONE",
        "StreamStatus": "ACTIVE",
        "StreamName": "Foo",
        "Shards": [
            {
                "ShardId": "shardId-000000000000",
                "HashKeyRange": {
                    "EndingHashKey": "340282366920938463463374607431768211455",
                    "StartingHashKey": "0"
                },
                "SequenceNumberRange": {
                    "StartingSequenceNumber": "49600945676829605166972608541704302084528788124637069314"
                }
            }
        ],
        "StreamARN": "arn:aws:kinesis:us-east-1:559483139917:stream/Foo",
        "EnhancedMonitoring": [
            {
                "ShardLevelMetrics": []
            }
        ],
        "StreamCreationTimestamp": 1572504543.0,
        "RetentionPeriodHours": 24
    }
}
```

Step4: check stream name

```
CHMB119868:Downloads samuelasirvat_r$ aws kinesis list-streams
[
    "StreamNames": [
        "Foo"
    ]
]
```

Step5: put the record using partitions

```
CHMB119868:Downloads samuelasirvat_r$ aws kinesis put-record --stream-name Foo --partition-key 123 --data testdata
{
    "ShardId": "shardId-000000000000",
    "SequenceNumber": "49600945676829605166972608706057767261137630198377742338"
}
```

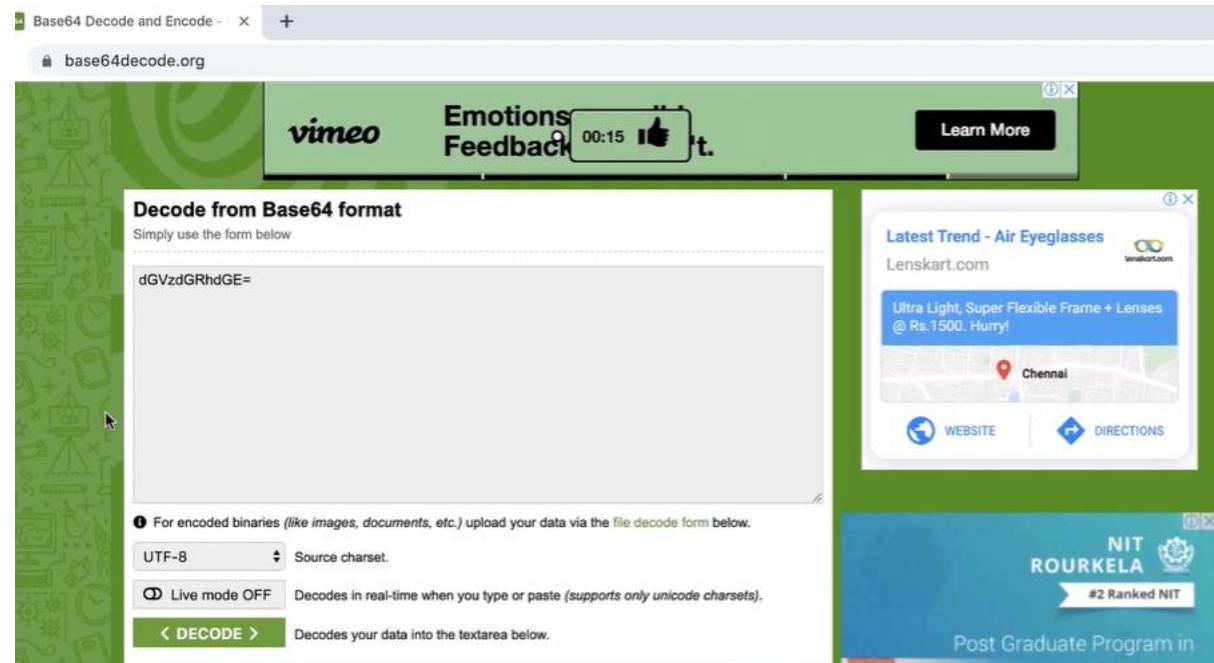
Step6:

```
CHMB119868:Downloads samuelasirvat_r$ aws kinesis get-shard-iterator --shard-id shardId-000000000000 --shard-iterator-type TRIM_HORIZON --stream-name Foo
{
    "ShardIterator": "AAAAAAAAAAH0zPVU4pRH2ocY5sj0YPbNntbU8eC9KtvpTZvCq9JTaKtzGC15jpadoJ2VCVxxmcvL4mHcBtjQNYlw0hPSRbSIyduYXe6nzFBwlX1c1LjQULdq3S+4uCIKkjwYUeX
j2609F9Y0NRJM15YDo39lRilHpVKZxif1LNQo0KvbLf1IH+apnC9u/qIBn+V7QN22876GzJ5g9kK1T4TJRtTEL4H"
```

Step7

```
CHMB119868:Downloads samuelasirvat_r$ aws kinesis get-records --shard-iterator AAAAAAAAHH0zPVU4pRH2ocY5sj0YPbNntbU8eC9KtvpTZvCq9JTaKtzGC15jpadoJ2VCVxxmcvL4mH
cBtjQNYlw0hPSRbSIyduYXe6nzFBwlX1c1LjQULdq3S+4uCIKkjwYUeXg2609F9Y0NRJM15YDo39lRilHpVKZxif1LNQo0KvbLf1IH+apnC9u/qIBn+V7QN22876GzJ5g9kK1T4TJRtTEL4H
{
    "Records": [
        {
            "Data": "dGVzdGRhdGE=",
            "PartitionKey": "123",
            "ApproximateArrivalTimestamp": 1572504625.14,
            "SequenceNumber": "49600945676829605166972608706057767261137630198377742338"
        }
    ],
    "NextShardIterator": "AAAAAAAAAAE8TeXj3p+QMpe8YEpZKv6/vGta9/dcSNUHlKR3VR3zVsqt7V73v+VyzuEBb+T38UnNjTDFsiyy4p0C1wGIiPiv856wabqCicIZtfQ2YCQD0g73unZQH7quM1m6
ehWeKoYzRAuzqij1Q6g83pdYztmODKZQ4fPpj18zkB9R5UhU2GK9dyGh6f361YHgnULBCUvpjj8YYo+R5EUoMjhjEe",
    "MillisBehindLatest": 0
}
```

Step8: decoding



Result: testdata (user enter the data)

Step 9: delete stream name

```
CHMB119868:Downloads samuelasirvat_r$ aws kinesis delete-stream --stream-name Foo  
CHMB119868:Downloads samuelasirvat_r$ █
```

Step9.1

```
CHMB119868:Downloads samuelasirvat_r$ aws kinesis describe-stream --stream-name Foo  
{  
    "StreamDescription": {  
        "KeyId": null,  
        "EncryptionType": "NONE",  
        "StreamStatus": "DELETING",  
        "StreamName": "Foo",  
        "Shards": [],  
        "StreamARN": "arn:aws:kinesis:us-east-1:559483139917:stream/Foo",  
        "EnhancedMonitoring": [  
            {  
                "ShardLevelMetrics": []  
            }  
        ],  
        "StreamCreationTimestamp": 1572504543.0,  
        "RetentionPeriodHours": 24  
    }  
}
```



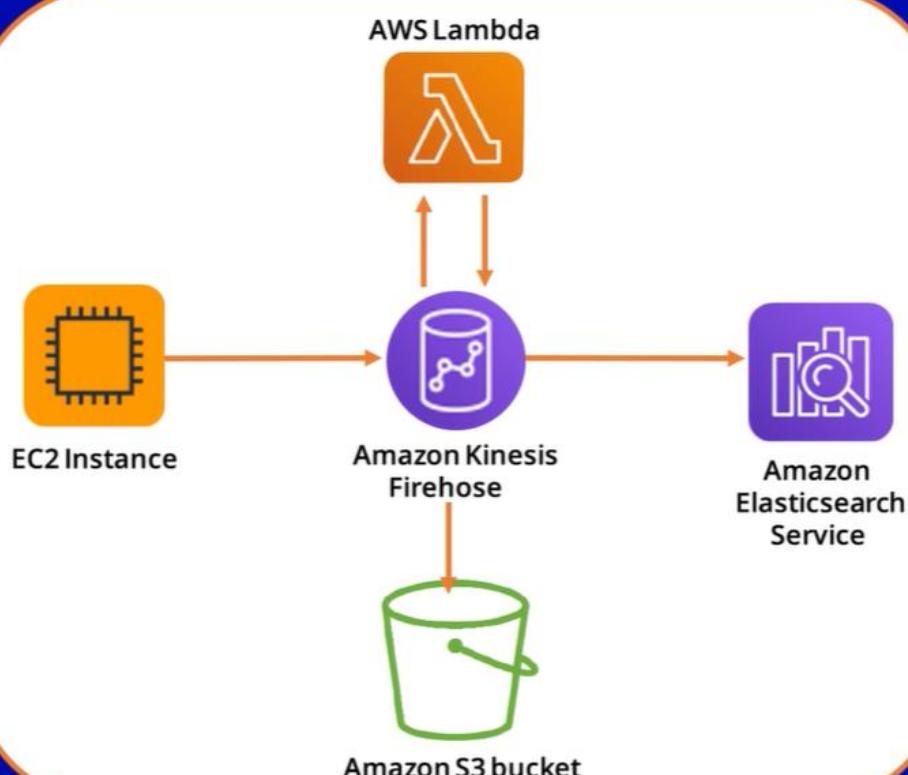
- Lambda example of service less computing which allows you run code without provision managing servers.
- Lambda you can run code anytype of app backend server all zero admission
- Just uploaded code lambda take care any things required to run scala code with high availability.
- You setup target automatically other aws services and call directly from any web or mobile app

Transferring Data using Lambda



Firehose can invoke a Lambda function to transform incoming data and then deliver the transformed data to S3, RedShift or Elasticsearch. And you can enable the firehose data transformation when you create your delivery stream.

Transferring Data using Lambda



Flow:

- Firehose incoming data buffering upto 3mb or buffering size specific delivery stream.
- Firehose involved specific lambda function each buffer batch sync.

- Transform data from lambda to firehouse for buffering
- Transfer data send to destination with specific buffering size and buffering interval reach.

Failure handling

- If data transformation fails unsuccessfully processed records are sent to your s3 bucket in processing failed folder.
- Fire hose retries invocation for three time and discard batch of records if retry is not successes
- Retry delivery period is 24 hours
- Invocation errors and records that have status of processing failed can be emitted to cloud watch logs.
- If data transfer fails to unsuccessful process record send to s3 bucket into process fail folder.

Data flow:

For s3 destination streaming records are delivered on your s3 bucket if data transformation is enabled, you could optionally backup supply data to any other s3 bucket.

Red shift destination:

Streaming records is kept in s3 bucket. To load data s3 bucket kinesis. Data issues an amazon red shift copy command into red shift cluster

In call transformation is enable you many optionally backup supply information every other Amazon s3 bucket.

Data delivery frequency in s3:

- In s3 buffer size (1mb to 128mb)
- Buffer interval 60 -900 sec
- Fire hose is capable of raising the buffers size dynamically to catch upto issues where data delivery to specific app is following behind.
- Retry delivery to your destination for upto 24 hrs.

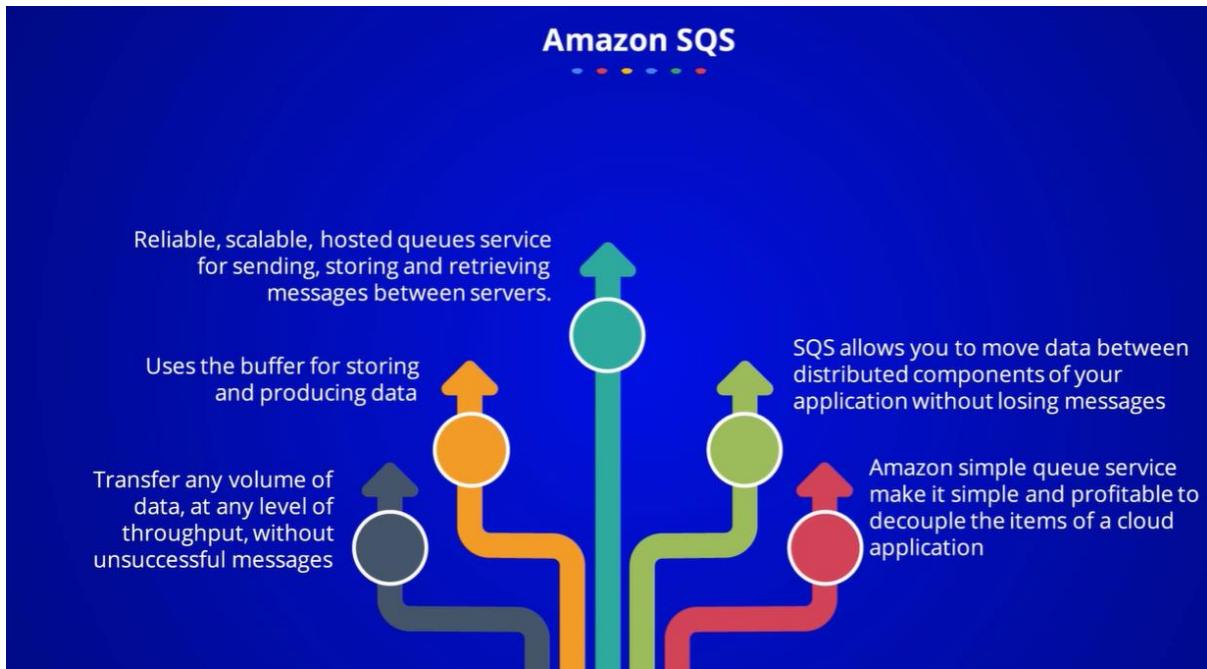
Data delivery frequency in Red Shift

- Red Shift frequency will depend on how fast red shift cluster finishes copy command.
- Retry in red shift 0 to 7200 sec.
- Fire hose will issues new copy cmd automatically when previous copy cmd has completed running.

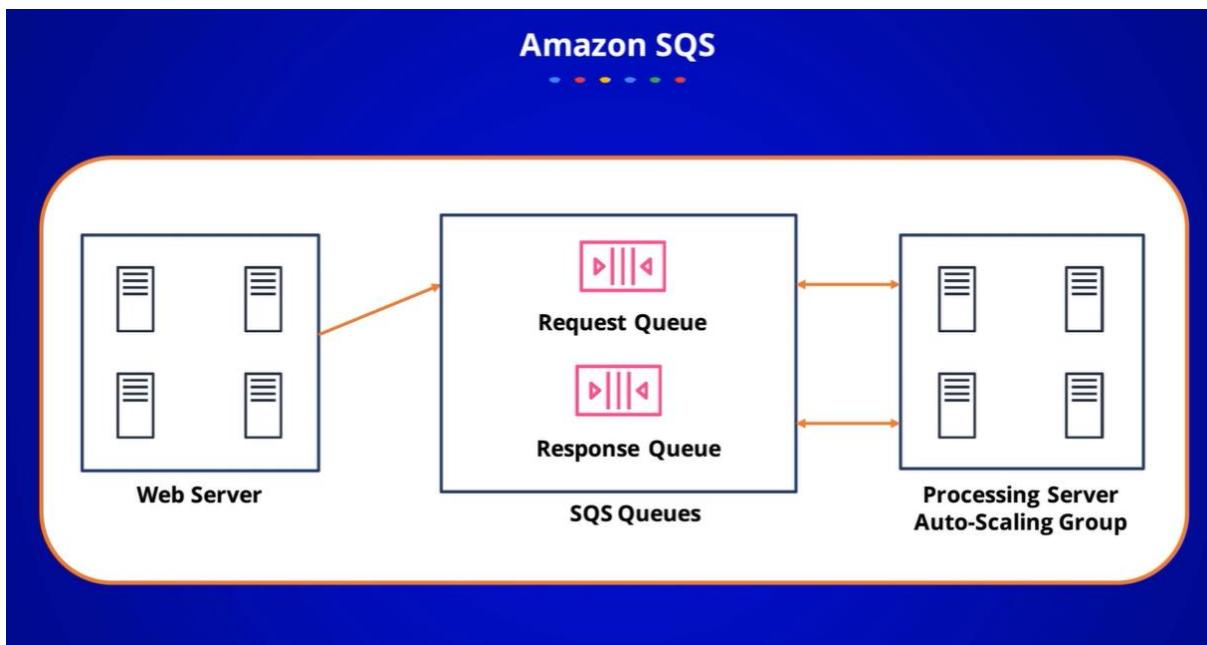
Data delivery frequency in elastic search

- Elastic search also depends on buffer size (1Mb to 100Mb) and buffer interval (60- 90 sec)
- Fire hose can raise size of buffer dynamically. If data delivery to elastic search is following behind.

- You can also specific retry duration from zero to 7200 sec.

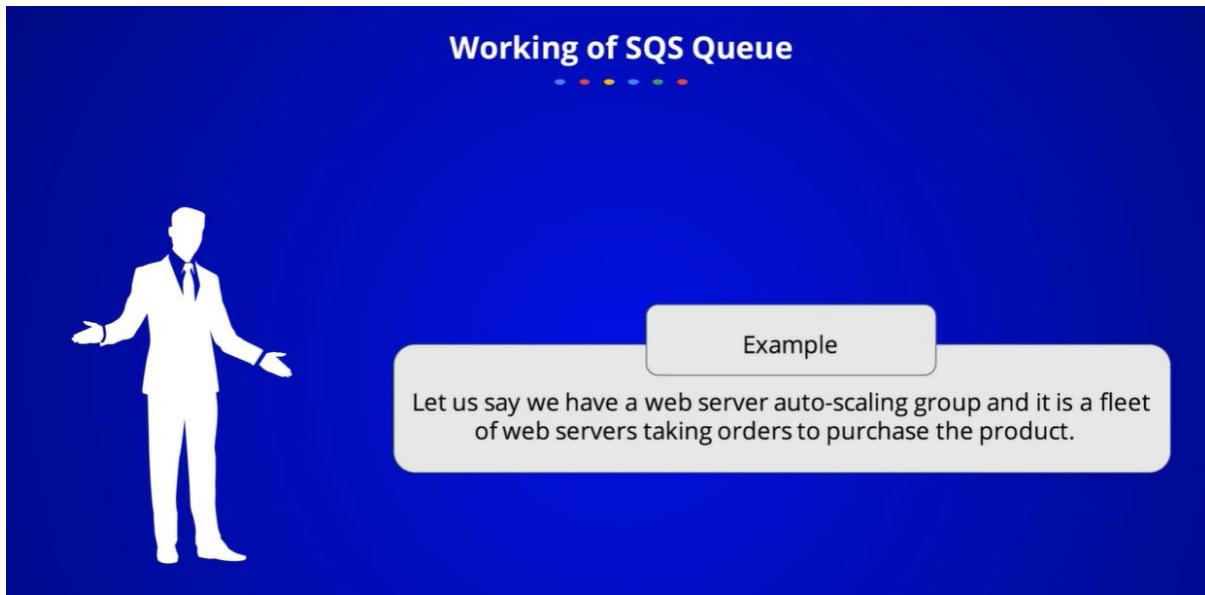


- Standards for simple queue services its reliable, scalable, post queue sending storing and retrieving message between server.

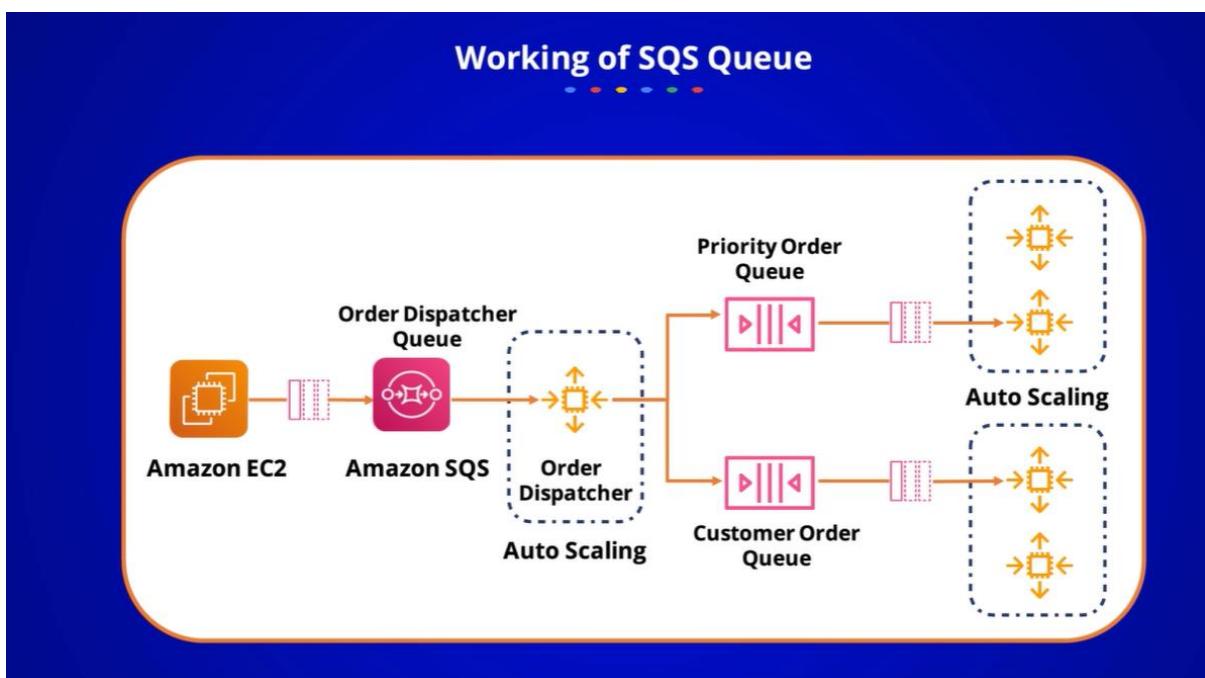


- Web server sub multiple reader/write which mean many complemt share in single queue.

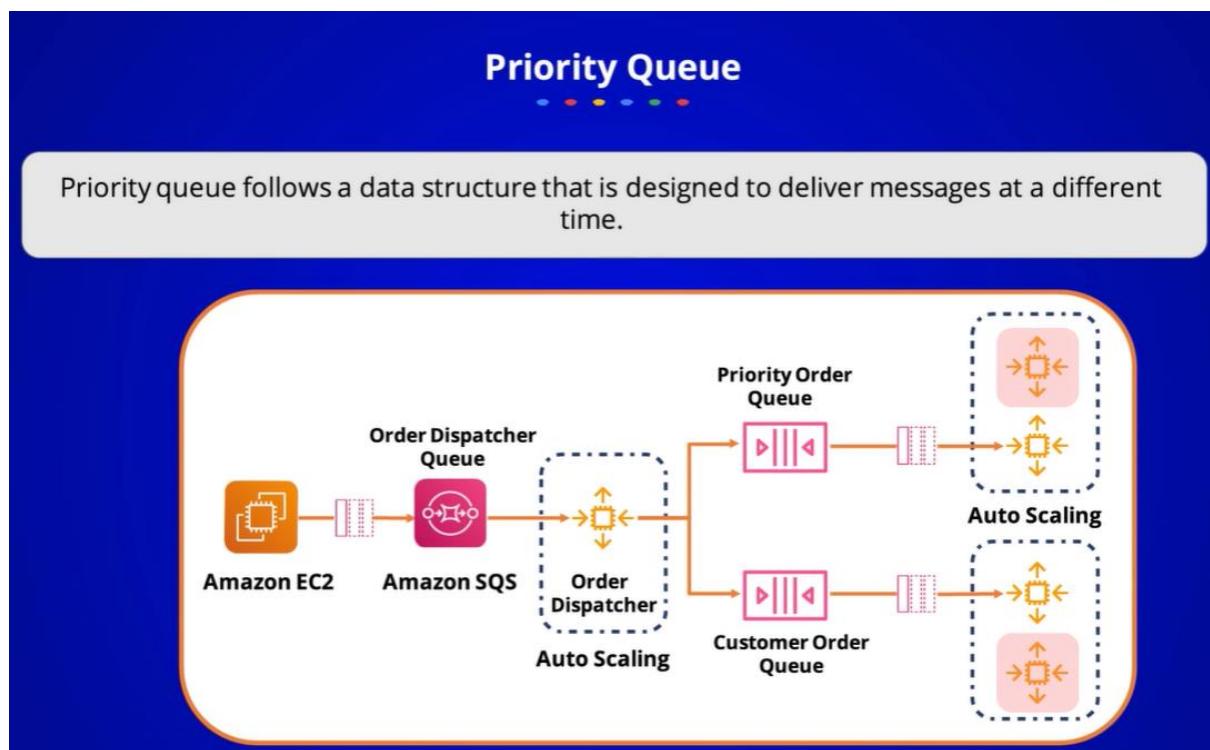
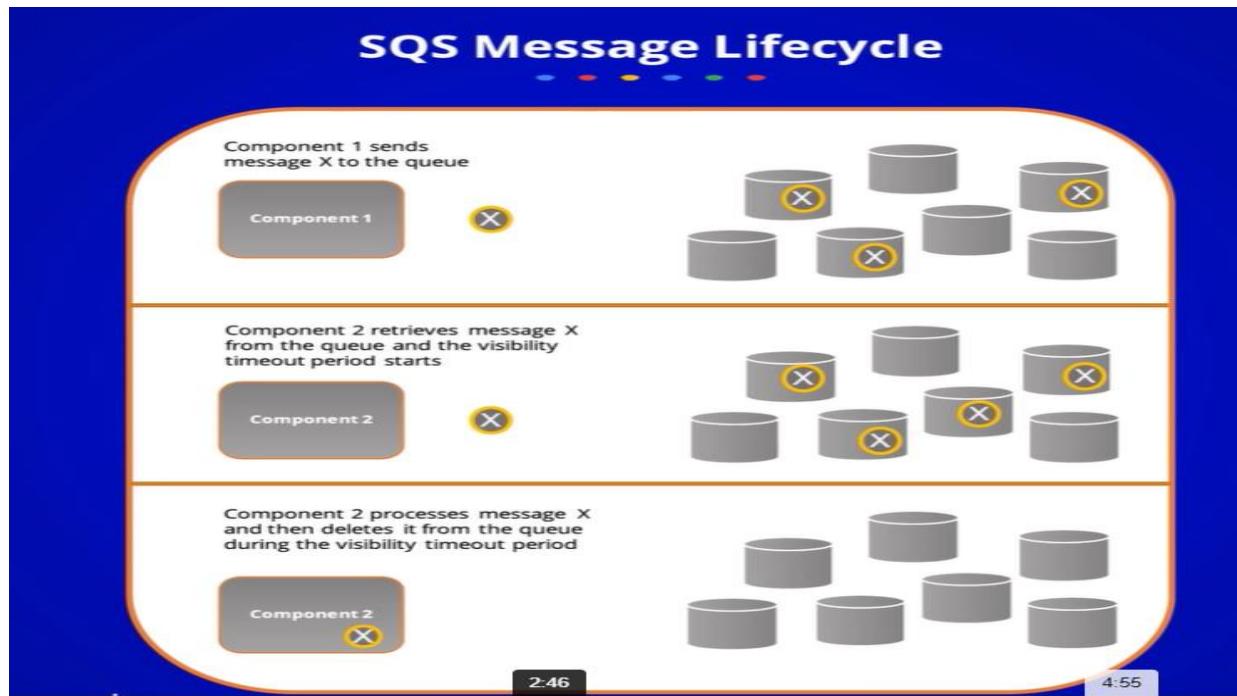
- Aws works fifo concept first in first out an fifo queue order send in message send once it remains available contain customer predict need that data
- Queue contain many region message can be contain queue upto 14 dayd this message send read simultaneously
- Sqs following long following its reduce minus cost while retrieve msg quickly and possible
- When queue is empty long fall request wait up to 20 sec next msg arrived.



- Web app place order into a sqs request queue this order set in queue until picked one free processing server. The server process server then back process to queue then customer.
- Queue can determine load on application. You can use length of collect determined queue and could watch determined load and could watch integrate with sqs queue to collect and view analyst metrics and you can several different discussion based on data.



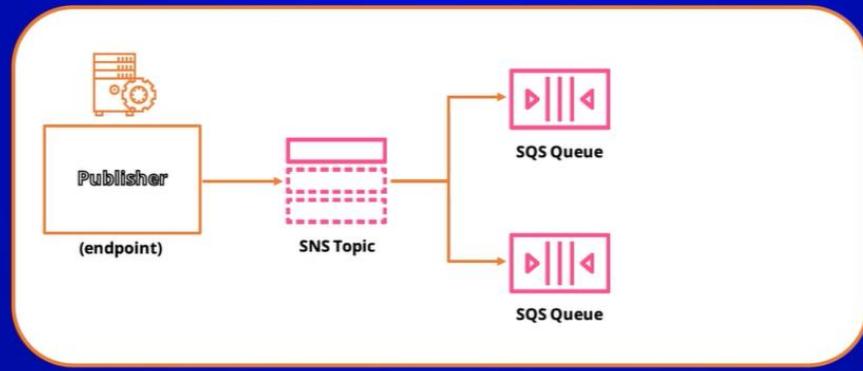
When web application place order sqs queue this order pickup one of the free processing server. The process order send back to priority order queue and send to customer order queue. The advantage of queue auto scaling scale back based on item sqs queue



Massage process two types high priority and low priority based on queue.

Priority Queue

The fan-out scenario involves SNS to send multiple topics using the SQS queue at the same time to decouple the processing further. It includes SNS messages which are sent to multiple topics, are replicated and pushed back in SQS queue to allow parallel processing.



AWS IoT



The AWS IoT service is a managed cloud platform that connects devices easily and securely to interact with cloud applications and other devices.

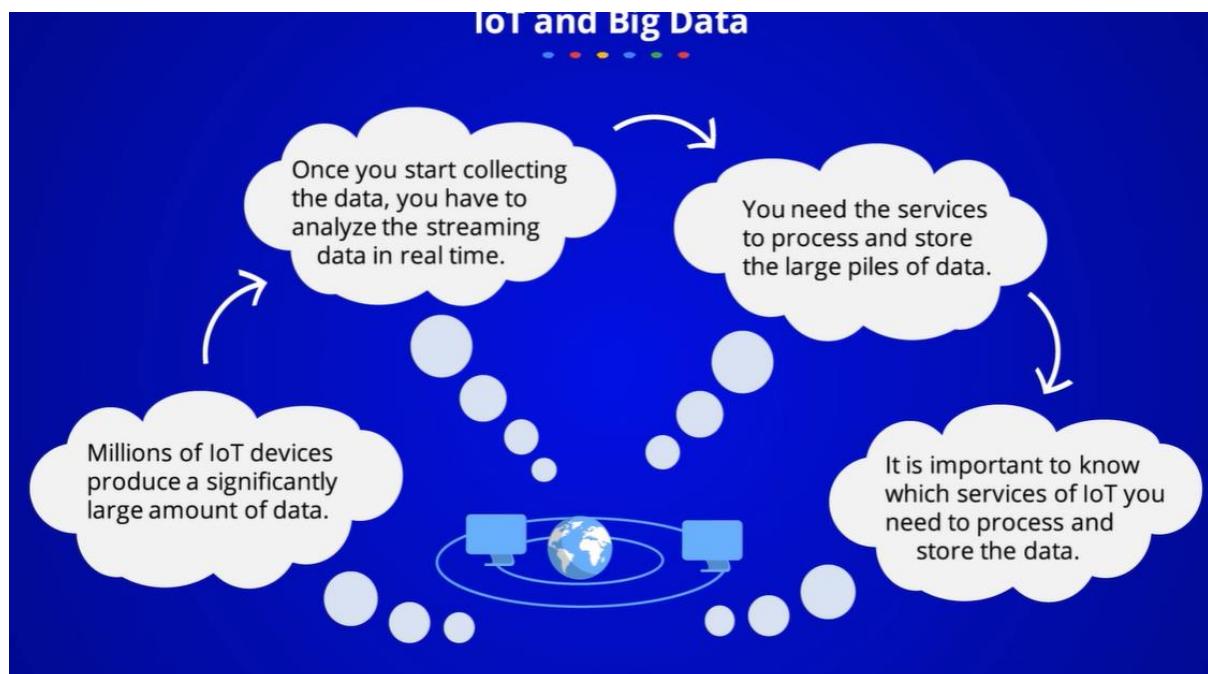
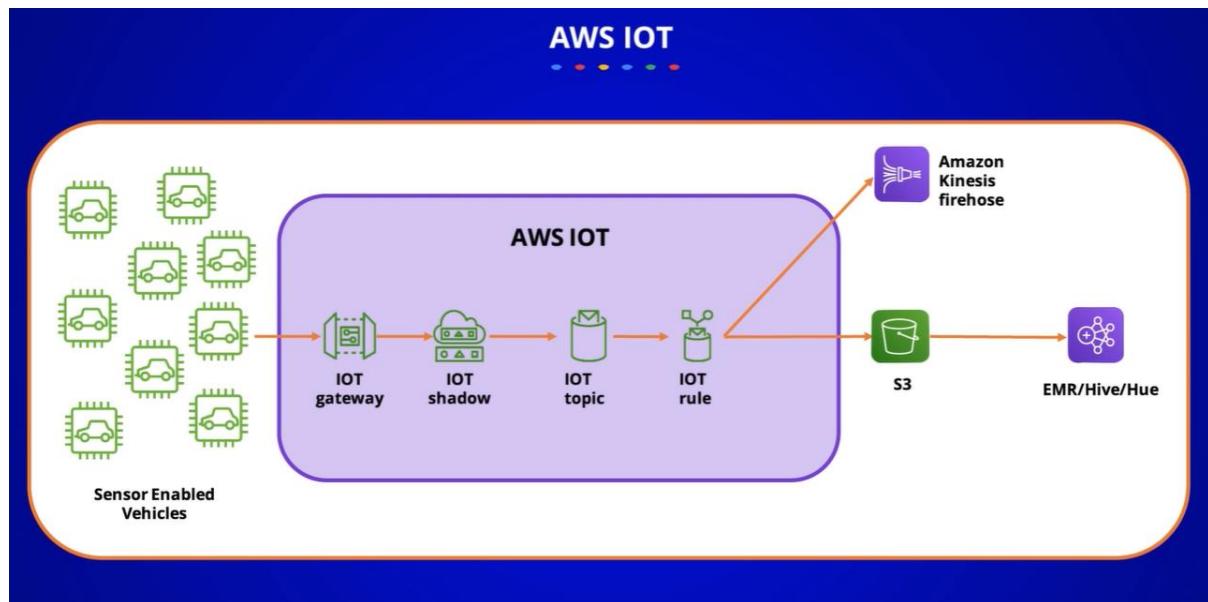
It supports billions of devices and trillions of messages and processes messages to other AWS services and other devices.

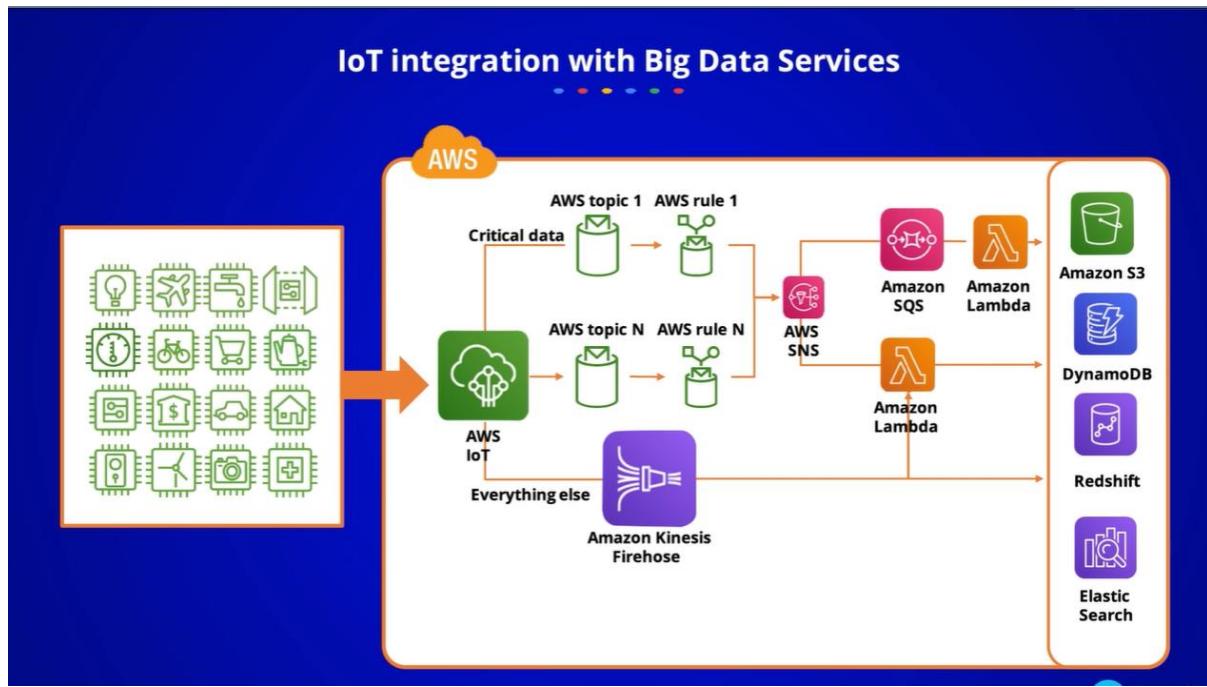


You can create applications that permit your customers to govern IoT devices from their phones or tablets.

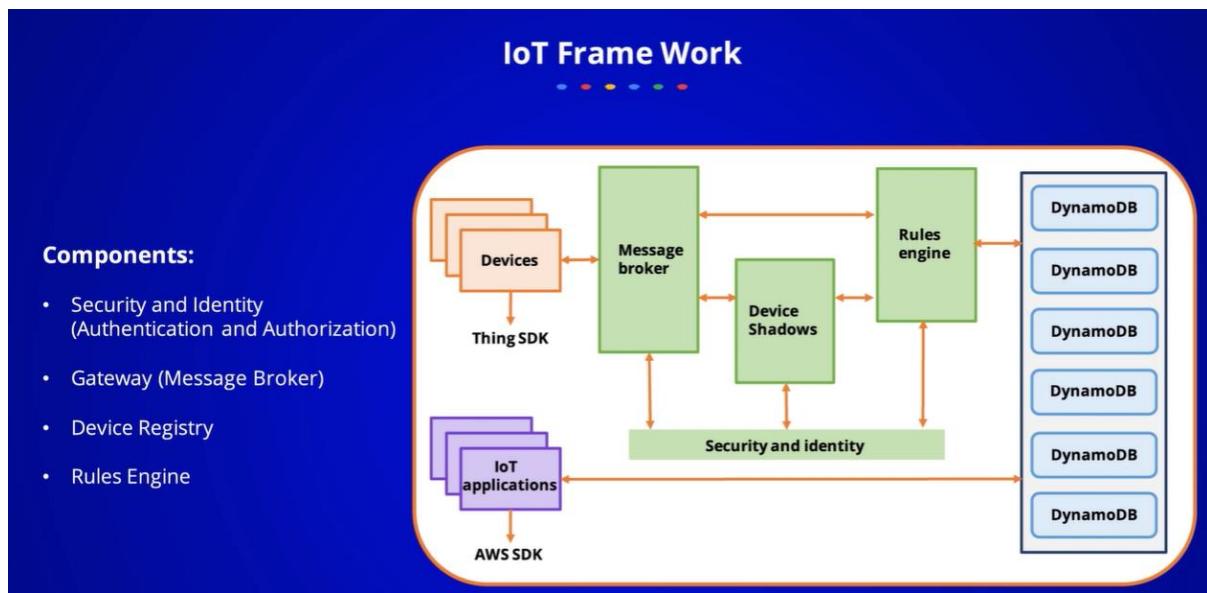


AWS IoT offers at ease, bi-directional communication among networked devices which include sensors, actuators, embedded microcontrollers, or smart home equipment and the AWS Cloud.





Once received message to IoT device. IoT device create rule and action data to elastic search service domain . data kinesis firehose stream write data into dynamo table even though Amazon machine learn base and production machine learning model. IoT machine model allow to change cloud watch alarm capture cloud watch metrics write data to s3 buckets, sqs queue sns push notification and finally invoked lambda function.

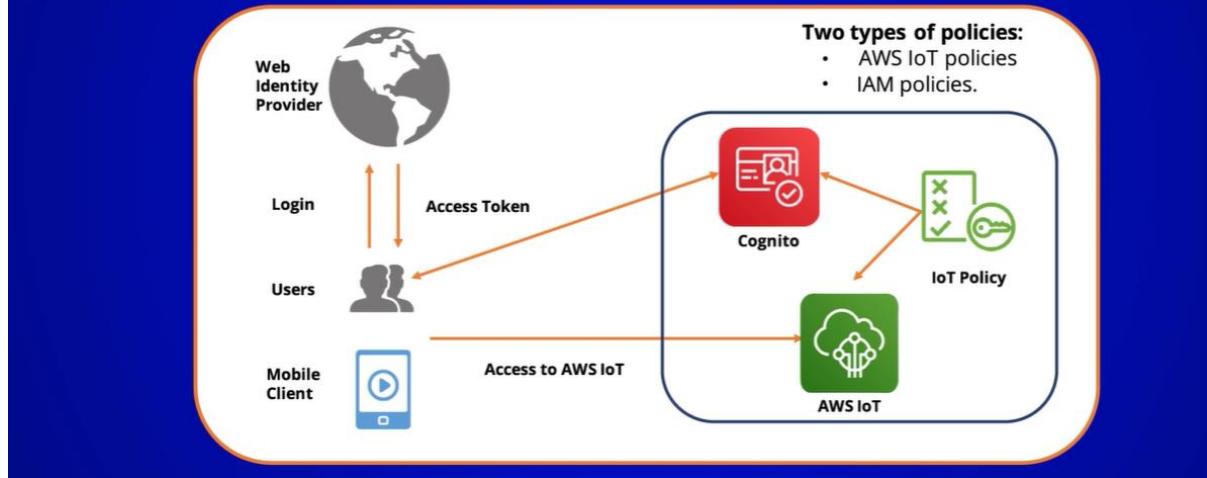


One device send to IoT device. IoT device having permission to write data into Elastic search data domain.

IoT Frame Work

Security and Identity (Authentication and Authorization)

Communication exchange between a tool and AWS IoT is included via the use of X.509 certificates.



Aws IoT cerificied and register and copy to our device. When device communicating aws IOT present certification identity credential.

Aws IoT use IAM policies use for user group and role. For mobile application authentication against IOT.

Authorisation two type of policy aws IoT policies and iam policies.

IoT Frame Work

.....

Gateway(Message Broker)

**It works to maintain the session and subscription
for all the connected devices in an IOT service.**



Gateway communicate one to one and one to many devices

The supported protocols are MQTT (Message Queue Telemetry Transport), Web sockets and HTTP

IoT Frame Work



Device Registry and Device Shadow

The device registry is a central location used for storing attributes

JSON document



Store and retrieve
the current state
for a thing

It acts as a message channel to send commands reliably to a thing
and stores the last known state in the AWS service.

IoT Frame Work



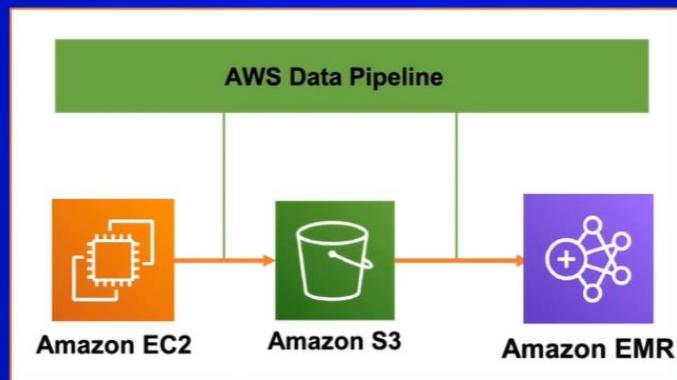
The Rule Engine collects the
statistics dispatched to the IoT cloud



Performs actions based totally on
factors which might be gathered data
and routes them to AWS endpoints

AWS Data Pipeline

A web provider that helps you in reliably processing and moving records between exceptional AWS compute and storage services as well as on-premises information sources at distinctive periods



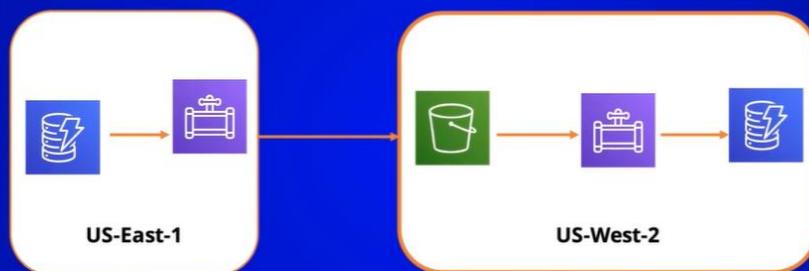
Data pipeline allows to create ETL work flow automate process moment of data schedule interval and terminate etl resource after work completed.

One of future data pipeline is ability to move data across region and you can copy entire dynamo db tables to another region or incremental copy table to another region.

You can reduce time frequency copy the table data pipeline help you reliable process move data between difference aws compute and store data services.

Example of a Data Pipeline

Consider an example, in which you have a DynamoDB table in us-east-1 and a data pipeline is used to export the data out of the table into an S3 bucket which resides in uswest-2 region.



You have another data pipeline in uswest-2 which imports that data into DynamoDB table from S3 bucket within the region.

Example of a Data Pipeline

The data pipeline is launching an EMR cluster to extract the data out of DynamoDB from us-east-1 and transfer it to S3 bucket in us-west-2.

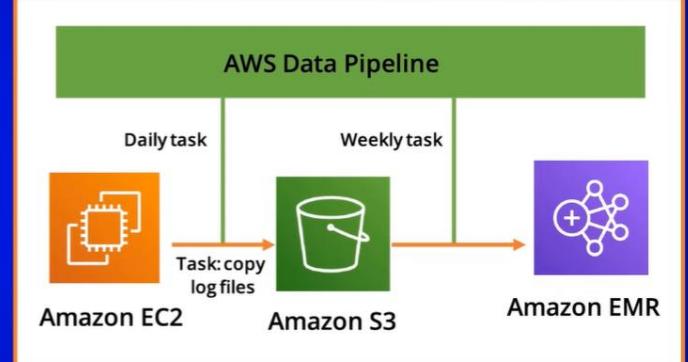


The data pipeline launches another EMR cluster to import the data into DynamoDB in the us-west-2 region.

Key Concept

Data Pipeline runs on an EC2 instance or an EMR cluster which are provisioned and terminated automatically.

The pipeline is the name of the container that contains things like data nodes, activities, preconditions, etc.



One of best benefits pipe line that can be scheduled all components move one location to another.

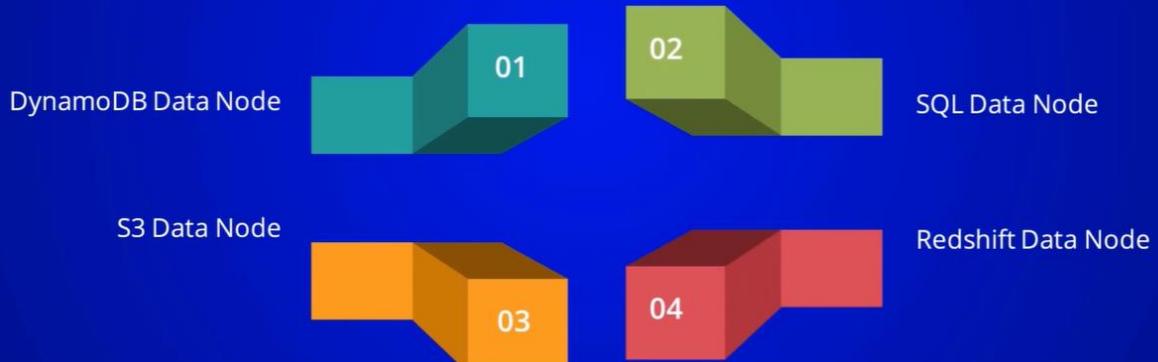
And pipe line run on premise run aws provide task run package install on own premise host.

As soon as install packet poll the pipe line works carrier out while time run select activity on premise source as example executing db store to process data base done. Aws pipe line issues to task run.

Data pipeline integrated with on premise environment.

Data Nodes

In AWS Data Pipeline, a data node defines the location and type of data that a pipeline activity uses as input or output. It is the end destination for your data.



dynamo db table contain data high activity or emr activity use.

Sql data note: sql databases query represent data for pipe line activity use.

Redshift data note: amazon redshift table contain redshift copy activity to use.

S3 data node: amzon location that contain one or more file contain for pipeline activity use.

Activity

Activity is an action that the data pipeline initiates on your behalf as part of the data pipeline. In the AWS Data Pipeline, an activity is a pipeline factor that defines the work to perform.



Copy Activity



Pig Activity



EMR Activity



Redshift Copy Activity



Hive Activity



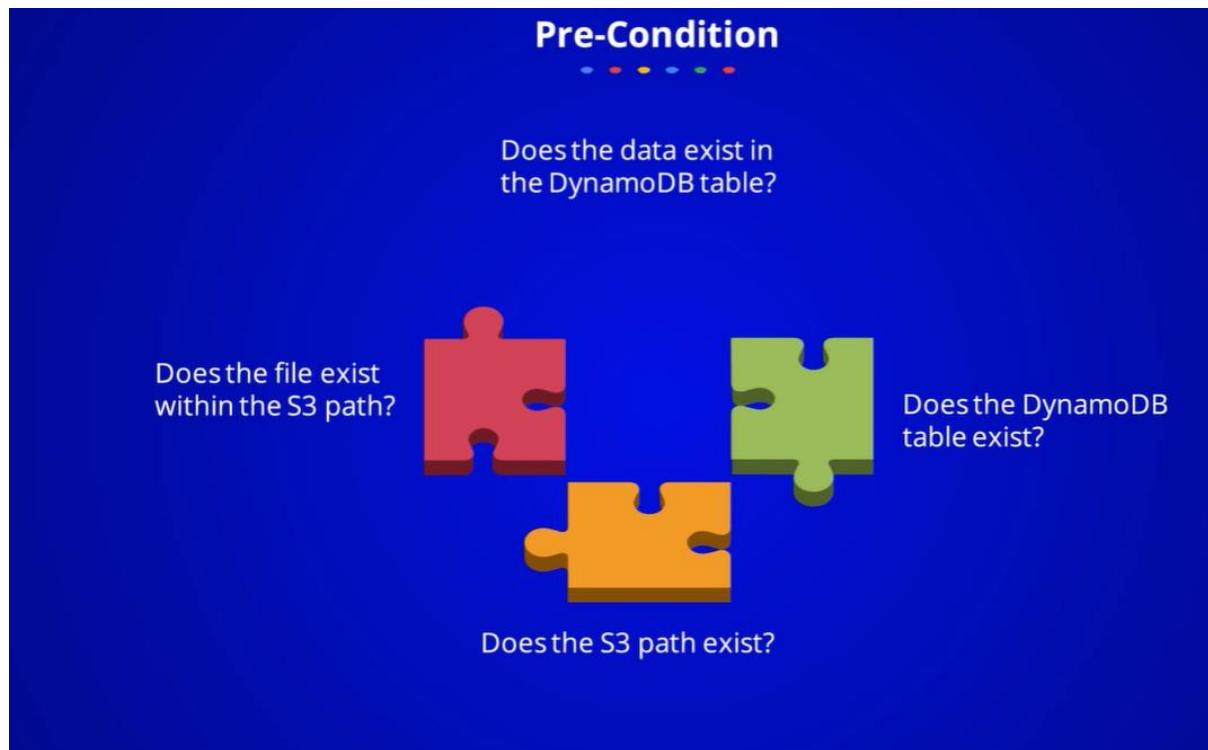
Shell Command Activity



Hive Copy Activity



SQL Activity



Screenshot of the AWS Lambda console showing the 'Create Pipeline' page. The top navigation bar includes 'Services', 'Resource Groups', and user information ('Samuel'). The main area shows the pipeline configuration:

Create Pipeline

You can create pipeline using a template or build one using the Architect page.

Name: S3-to-dynamodb-copy

Description (optional): (empty)

Source: Build using a template
Import DynamoDB backup data from S3

Import a defi... Loading template...
 Build using Architect

Schedule

You can run your pipeline once or specify a schedule. More

Run: on pipeline activation on a schedule

Run every: 1 day(s)

Starting: on pipeline activation
2019-10-30 13:49 UTC (Current time is 13:50 UTC)
YYYY-MM-DD HH:MM

Ending: never
 after 1 occurrence(s)
2019-10-31 13:49 UTC (Current time is 13:50 UTC)

e **Build using a template**

✓ Choose...
 Getting Started
 Getting Started using ShellCommandActivity
 AWS Command Line Interface (CLI) Templates
 Run AWS CLI command
 DynamoDB Templates
 Export DynamoDB table to S3
Import DynamoDB backup(.data) from S3

sche
 Elastic MapReduce (EMR) Templates
 Run job on an Elastic MapReduce cluster

n
 RDS Templates
 Full copy of RDS MySQL table to S3
 Incremental copy of RDS MySQL table to S3

y
 Load S3 data into RDS MySQL table
 Redshift Templates

g
 Full copy of RDS MySQL table to Redshift
 Incremental copy of RDS MySQL table to Redshift
 Load data from S3 into Redshift

Schedule

on a schedule

Run every day(s)

Starting on pipeline activation
 YYYY-MM-DD HH:MM UTC (Current time is 13:50 UTC)

Ending never
 after occurrence(s)
 YYYY-MM-DD HH:MM UTC (Current time is 13:50 UTC)

Pipeline Configuration

Logging Enabled
 Disabled

Copy execution logs to S3. [More](#)

S3 location for logs

Security/Access

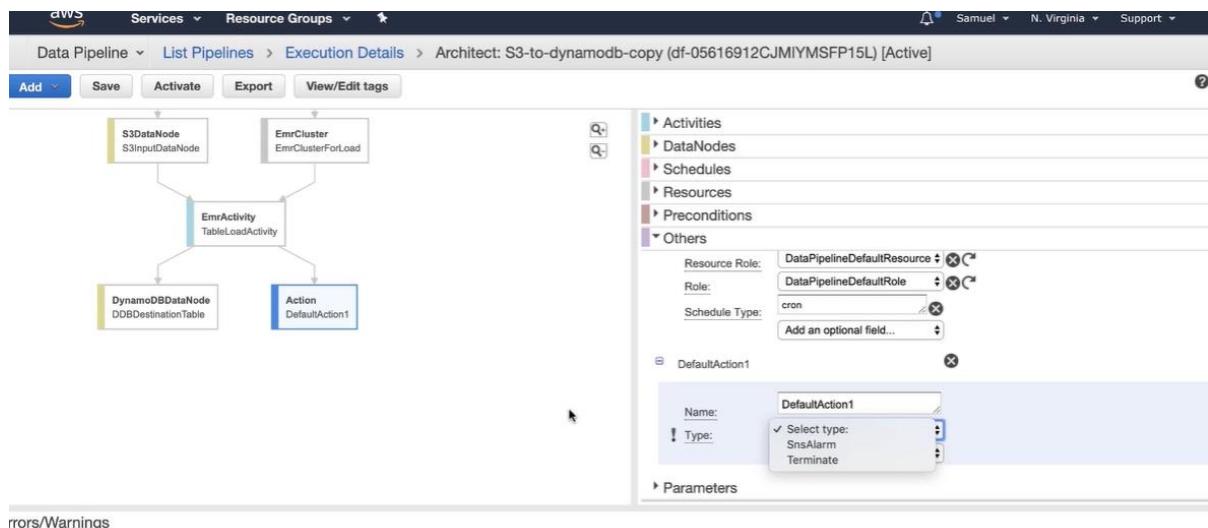
IAM roles Default
 Custom

IAM Roles let you control permissions for AWS Data Pipeline and your EC2 applications. [More](#)

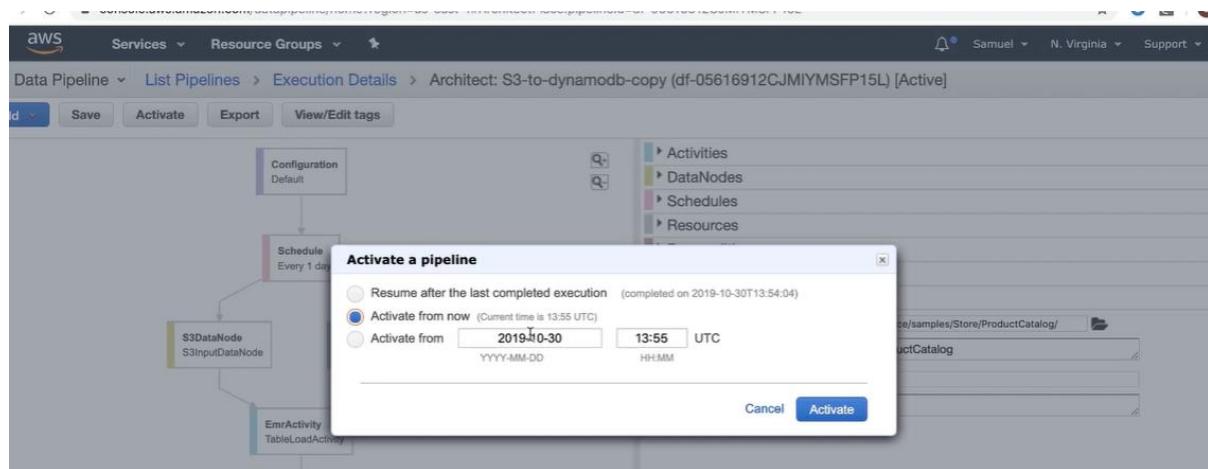
Tags

i Add up to 10 tags to your pipeline. These tags will be applied to the pipeline as well as any resources created by the pipeline. A tag consists of a case-sensitive key-value pair. [Learn more](#)

Key	Value (Optional)



Errors/Warnings



Create DynamoDB table

[Tutorial](#) [?](#)

DynamoDB is a schema-less database that only requires a table name and primary key. The table's primary key is made up of one or two attributes that uniquely identify items, partition the data, and sort data within each partition.

Table name*	<input type="text" value=""/>	i
Primary key*	Partition key	<input type="text"/> String i
<input type="checkbox"/> Add sort key		

Table settings

Default settings provide the fastest way to get started with your table. You can modify these default settings now or after your table has been created.

Use default settings

- No secondary indexes.
- Auto Scaling capacity set to 70% target utilization, at minimum capacity of 5 reads and 5 writes.
- Encryption at Rest with DEFAULT encryption type.

[+ Add tags NEW!](#)

Additional charges may apply if you exceed the AWS Free Tier levels for CloudWatch or Simple Notification Service. Advanced alarm settings are available in the CloudWatch management console.

[Cancel](#) [Create](#)

console.aws.amazon.com/dynamodb/name/region=us-east-1#tables:selected=copyfroms3-ProductCatalog;tab=items

Services Resource Groups * Samuel N. Virginia Support

[Create table](#) [Delete table](#)

Filter by table name Choose a table ... Actions

Name

copyfroms3-ProductCatalog

copyfroms3-ProductCatalog Close

Overview Items Metrics Alarms Capacity Indexes Global Tables Backups Triggers More

Scan: [Table] copyfroms3-ProductCatalog: Id ▾ Viewing 0 to 0 items

Scan [Table] copyfroms3-ProductCatalog: Id Add filter Start search

Id [i](#)

Services Resource Groups * Samuel N. Virginia Support

[Create table](#) [Delete table](#)

Filter by table name Choose a table ... Actions

Name

copyfroms3-ProductCatalog

copyfroms3-ProductCatalog Close

Overview Items Metrics Alarms Capacity Indexes Global Tables Backups Triggers More

Scan: [Table] copyfroms3-ProductCatalog: Id ▾ Viewing 1 to 8 items

Scan [Table] copyfroms3-ProductCatalog: Id Add filter Start search

Id i	Price	ProductCategory	Title	BicycleType	Brand	Color
101	2	Book	Book 101 Title			
102	20	Book	Book 102 Title			
103	2000	Book	Book 103 Title			
201	100	Bicycle	18-Bike-201 i	Road	Mountain A	{ "Black"
202	200	Bicycle	21-Bike-202	Road	Brand-Company A	{ "Black"
203	300	Bicycle	19-Bike-203	Road	Brand-Company B	{ "Black"
204	400	Bicycle	18-Bike-204	Mountain	Brand-Company B	{ "Red"
205	500	Bicycle	20-Bike-205	Hybrid	Brand-Company C	{ "Black"



Lesson 4:

AWS Big Data Storage Services



Amazon Glacier



Exceptionally low-cost service of AWS which offers secure and substantial storage for data archiving and backup.



It is optimized to hold expenses low for the data which is accessed infrequently, and for which retrieval time of numerous hours is appropriate.

Amazon Glacier



Current archiving Problem in Traditional storage

- Companies generally over-pay for archival of data.
- Provision for an unexpected increase

Amazon Glacier game Changer

- You do not pay anything upfront
- Pay a very low fee for storage
- Scale your utilization up or down as needed

Glacier and Big Data



Big data uses data operating at petabyte scale

1

Data needs to be archived at a lower cost

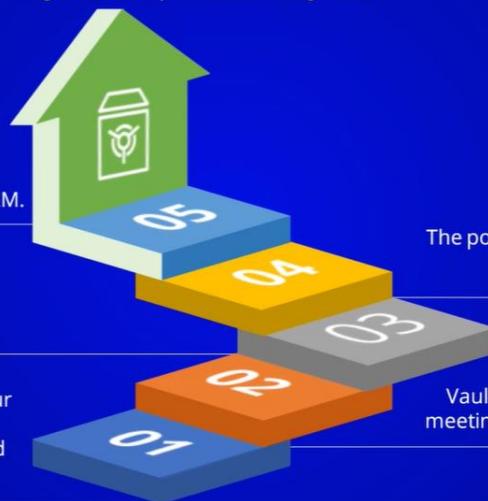
2

Glacier will provide you a greater benefit in terms of objects and using object life cycle policy in S3

3

Glacier and Big Data

Bigdata Compliance through Glacier:



The policies are created by using IAM.

Vault Lock policy used to deploy and enforce control.

Compliance requirement: Keep your data for a specific amount of time and during this time, the data could be read but not changed.

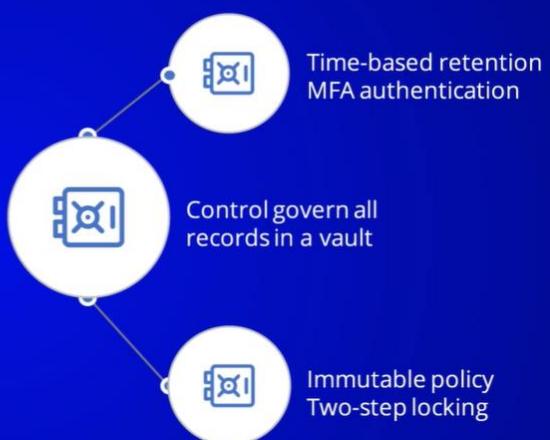
The policies can be locked from editing.

Vault Lock feature in Glacier helps in meeting the compliance requirements.

Glacier and Big Data

Compliance Storage with Vault Lock

Amazon Glacier Vault Lock gives you the ease of setting compliance controls on individual vaults and enforce them via a lockable policy.



Glacier and Big Data

Management features: Vault Lock



- Non-overwrite, non-erasable records
- Time-based retention with "AchieveAgeInDays" control
- Policy lockdown (strong governance)
- Legal hold with vault-level tags
- Optionally configure third-party access and provides temporary access



EXAM TIP: The Vault Policies are implemented by using IAM.

Glacier and Big Data

Accessing Amazon Glacier

Direct Amazon Glacier API/SDK Amazon S3 lifecycle integration Third-party tools and gateways



DynamoDB Introduction



DynamoDB is fast, flexible, and fully managed NoSQL database service for the applications that need consistent and scalable latency.

DynamoDB



Document Data Model



Key-Value Data Model

DynamoDB Introduction



Mobile and web application



Gaming



Internet of Things



Live online voting

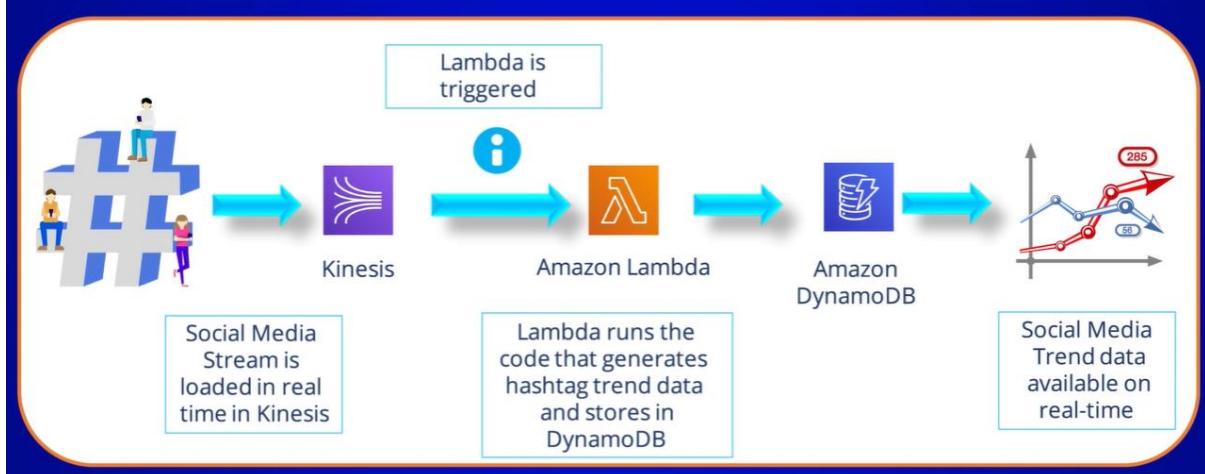


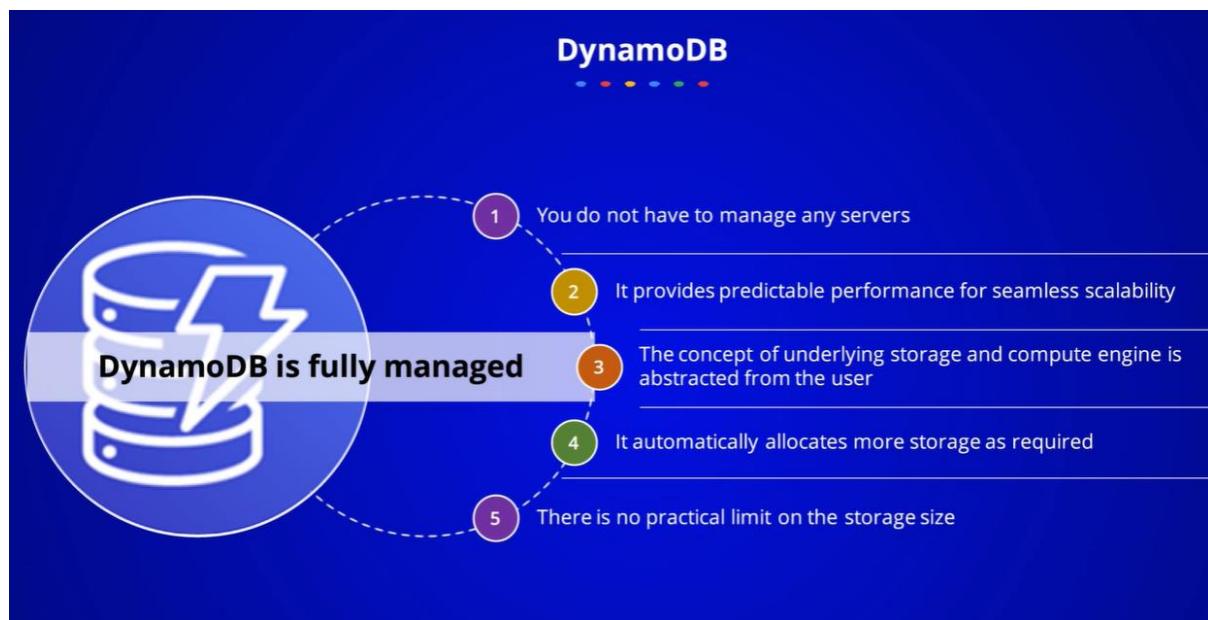
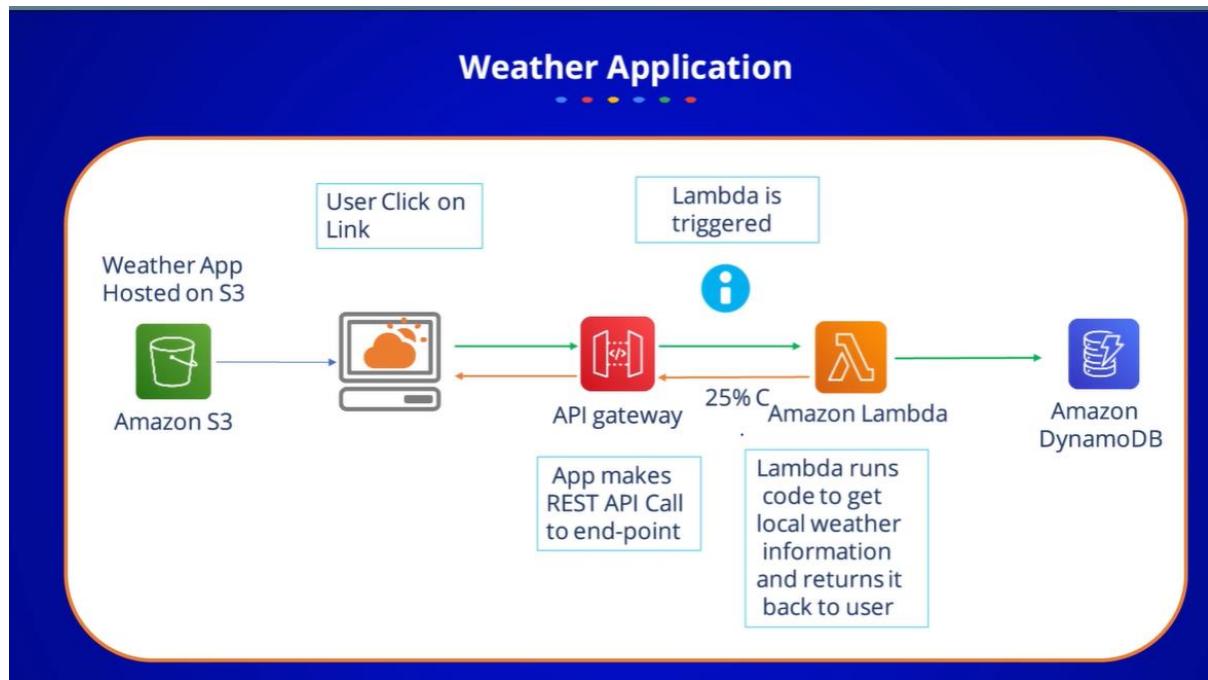
Session management

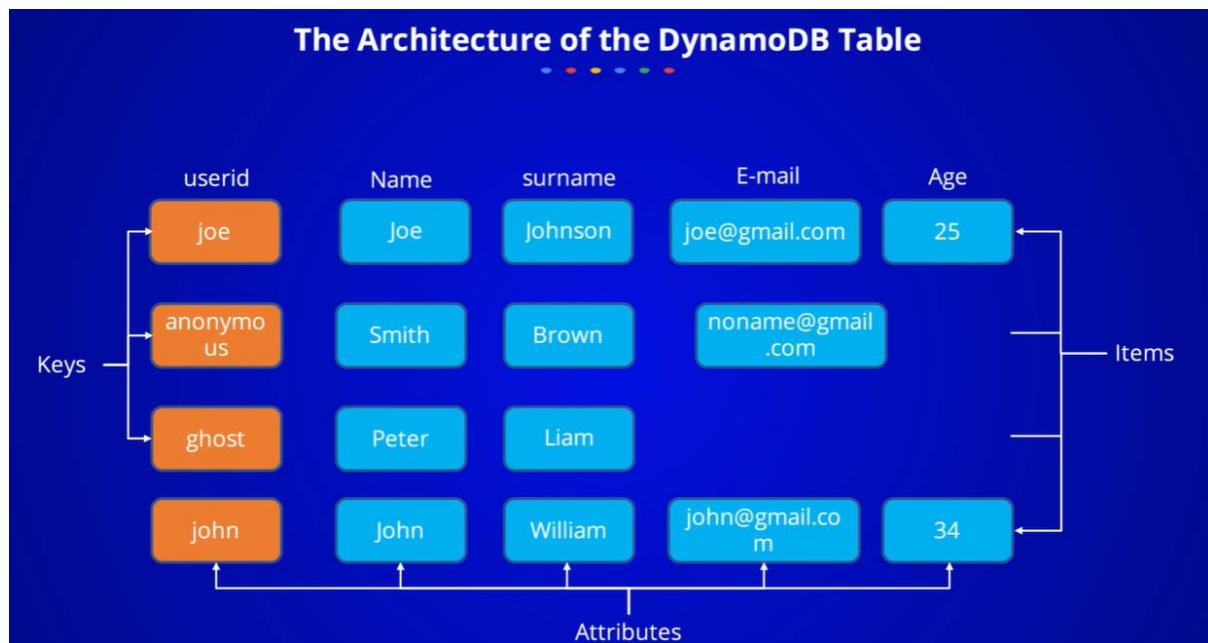
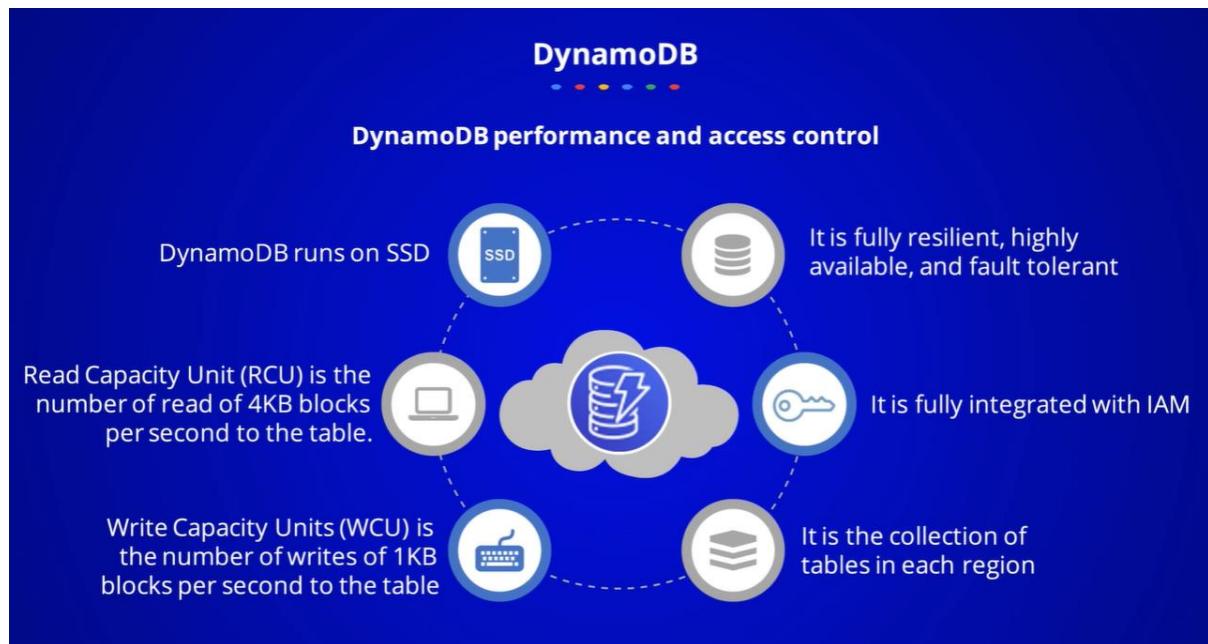


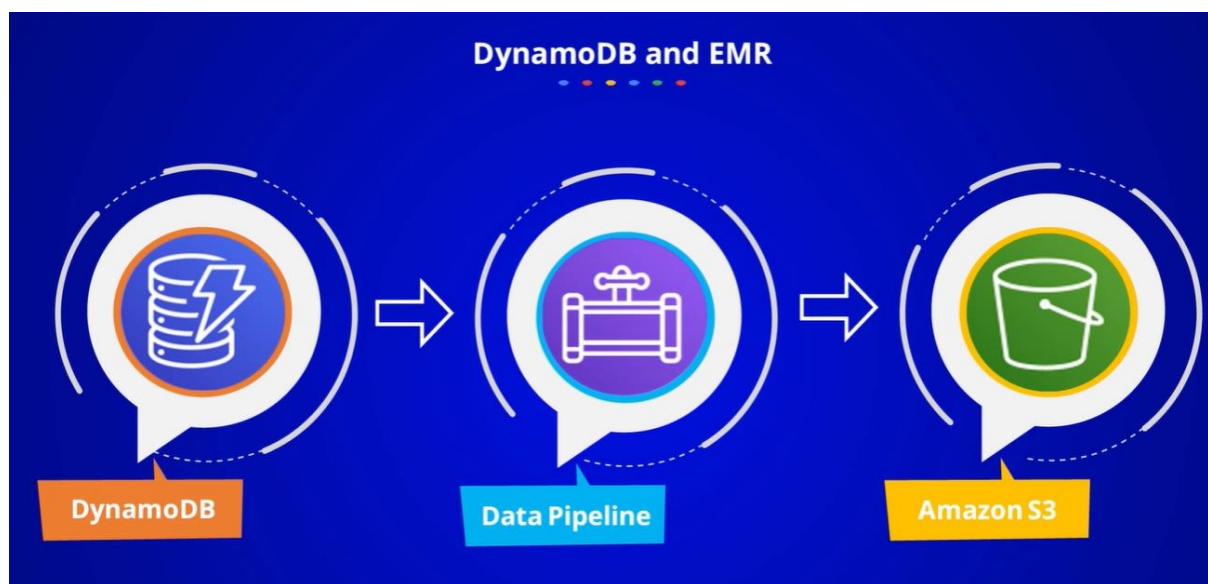
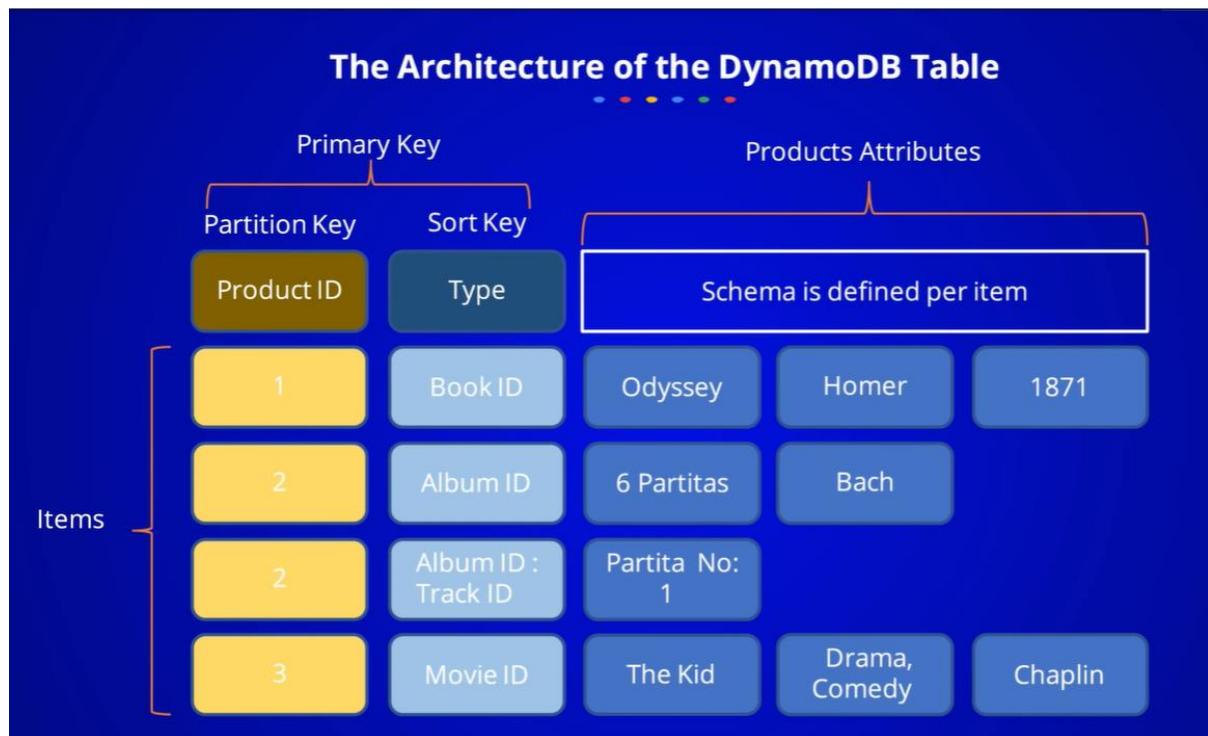
Storing S3 object metadata in the DynamoDB table.

Analysis of Streaming Social Media Data









Using copy cmd directly copy from dynamo db to redshift.

Dynamodb export data using data pipeline emr cluster, emr running dynamo db table write into s3 bucket. If s3 want write data to dynamo db using data pipeline emr cluster write with dynamo db.