

SAA-6

AWS Integration & Messaging

Section Introduction

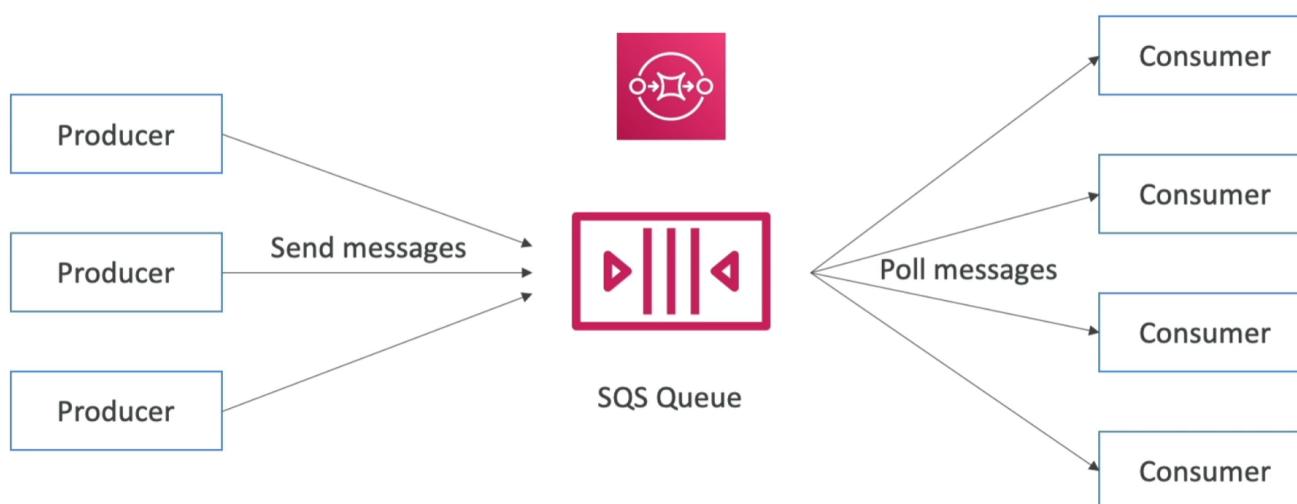
- When we start deploying multiple applications, they will inevitably need to communicate with one another
- There are two patterns of application communication



- Synchronous between applications can be problematic if there are sudden spikes of traffic
- What if you need to suddenly encode 1000 videos but usually it's 10?
- In that case, it's better to **decouple** your applications,
 - using SQS: queue model
 - using SNS: pub/sub model
 - using Kinesis: real-time streaming model
- These services can scale independently from our application

Amazon SQS

What's a queue?



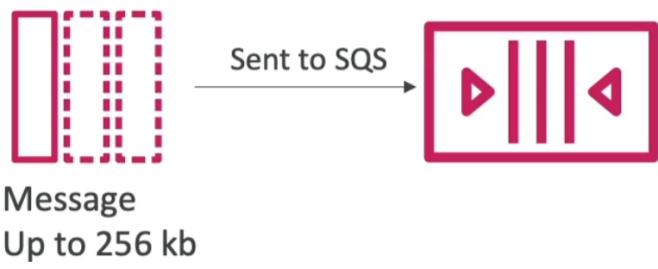
Amazon SQS - Standard Queue

- Oldest offering (over 10 years old)
- Fully managed service, used to **decouple applications**
- Attributes:
 - Unlimited throughput, unlimited number of messages in queue
 - Default retention of messages: 4 days, maximum of 14 days
 - Low latency (< 10 ms on publish and receive)
 - Limitation of 256KB per message sent
- Can have duplicate messages (at least once delivery, occasionally)
- Can have out of order messages (best effort ordering)

SQS - Producing Messages

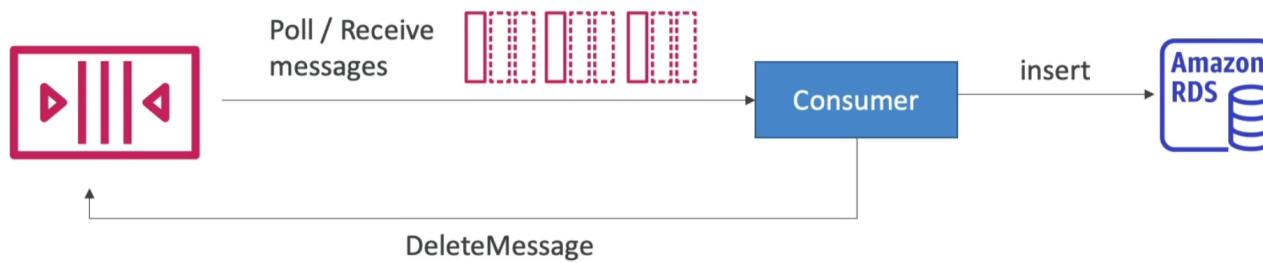
- Produced to SQS using the SDK (SendMessage API)
- The message is **persisted** in SQS until a consumer deletes it
- Message retention: default 4 days, up to 14 days
- Example: send an order to be processed

- Order id
- Customer id
- Any attributes you want
- SQS standard: unlimited throughput



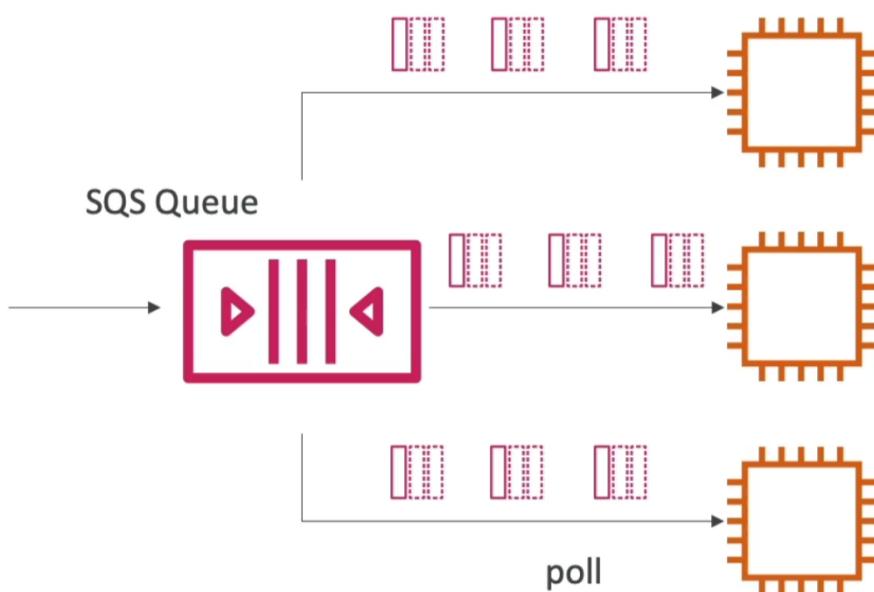
SQS - Consuming Messages

- Consumers (running on EC2 instances, servers, or AWS Lambda)...
- Poll SQS for messages (receive up to 10 messages at a time)
- Process the messages (example: insert the message into an RDS database)
- Delete the messages using the DeleteMessage API

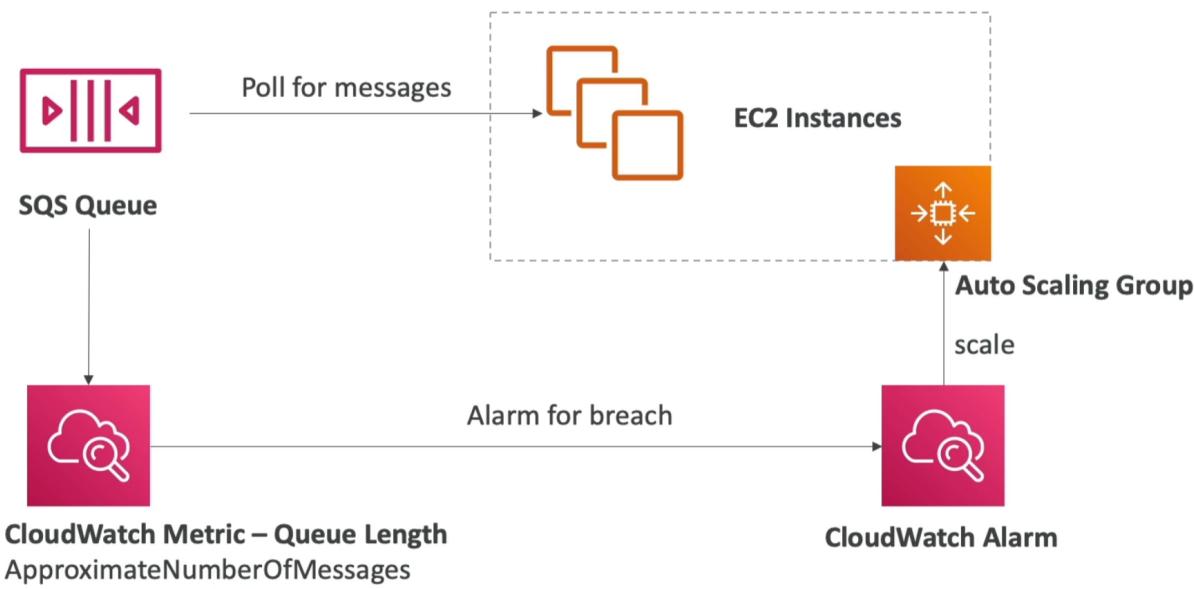


SQS - Multiple EC2 Instances Consumers

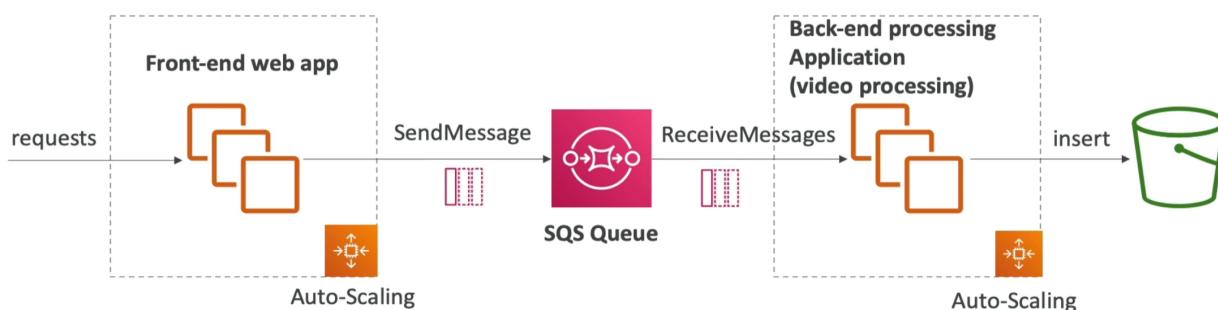
- Consumers receive and process messages in parallel
- At least once delivery
- Best-effort message ordering
- Consumers delete messages after processing them
- We can scale consumers horizontally to improve throughput of processing



SQS with Auto Scaling Group (ASG)



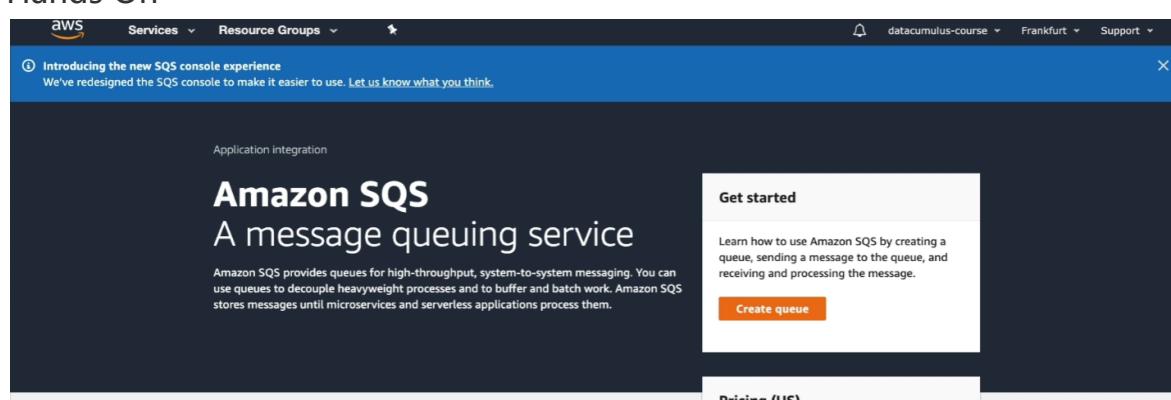
SQS to decouple between application tiers



Amazon SQS - Security

- **Encryption:**
 - In-flight encryption using HTTPS API
 - At-rest encryption using KMS keys
 - Client-side encryption if the client wants to perform encryption/decryption itself
- **Access Controls:** IAM policies to regulate access to the SQS API
- **SQS Access Policies** (similar to S3 bucket policies)
 - Useful for cross-account access to SQS queues
 - Useful for allowing other services (SNS, S3 ...) to write to an SQS queue

Hands On



Create queue

Details

Type

Choose the queue type for your application or cloud infrastructure.

You can't change the queue type after you create a queue.

Standard [Info](#)

At-least-once delivery, message ordering isn't preserved

- At-least once delivery
- Best-effort ordering

FIFO [Info](#)

First-in-first-out delivery, message ordering is preserved

- First-in-first-out delivery
- Exactly-once processing

Name

MyQueue

A queue name is case-sensitive and can have up to 80 characters. You can use alphanumeric characters, hyphens (-), and underscores (_).

Configuration

Set the maximum message size, visibility to other consumers, and message retention. [Info](#)

Visibility timeout [Info](#)

30

Message retention period [Info](#)

4

Should be between 0 seconds and 12 hours.

Should be between 1 minute and 14 days.

Delivery delay [Info](#)

0

Maximum message size [Info](#)

256

Should be between 0 seconds and 15 minutes.

Should be between 1 KB and 256 KB.

Receive message wait time [Info](#)

0

Should be between 0 and 20 seconds.

Access policy

Define who can access your queue. [Info](#)

Choose method

Basic

Use simple criteria to define a basic access policy.

Advanced

Use a JSON object to define an advanced access policy.

Define who can send messages to the queue

Only the queue owner

Only the owner of the queue can send messages to the queue.

Only the specified AWS accounts, IAM users and roles

Only the specified AWS account IDs, IAM users and roles can send messages to the queue.

Define who can receive messages from the queue

Only the queue owner

Only the owner of the queue can receive messages from the queue.

Only the specified AWS accounts, IAM users and roles

Only the specified AWS account IDs, IAM users and roles can receive messages from the queue.

JSON (read-only)

```
{
  "Version": "2008-10-17",
  "Id": "__default_policy_ID",
  "Statement": [
    {
      "Sid": "__owner_statement",
      "Effect": "Allow",
      "Principal": {
        "AWS": "387124123361"
      },
      "Action": [
        "SQS:*"
      ],
      "Resource": "arn:aws:sqs:eu-central-1:387124123361:DemoQueue"
    }
  ]
}
```

▼ Encryption - Optional

Amazon SQS provides in-transit encryption by default. To add at-rest encryption to your queue, enable server-side encryption. [Info](#)

Server-side encryption

Disabled

Enabled

Choose CMK

Choose a CMK alias

Enter the CMK alias

Customer master key [Info](#)

(Default) alias/aws/sqs

Description

Default master key that protects my SQS messages when no other key is defined

Account

387124123361

CMK ARN

arn:aws:kms:eu-central-1:387124123361:key/c7582c86-f57a-42c5-8c91-ca9617d7fa56

Data key reuse period

5

Should be between 1 minute and 24 hours.

► **Dead-letter queue - Optional**
Send undeliverable messages to a dead-letter queue. [Info](#)

► **Tags - Optional**
A tag is a label assigned to an AWS resource. Use tags to search and filter your resources or track your AWS costs. [Learn more](#)

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

DemoQueue

[Edit](#) [Delete](#) [Purge](#) [Send and receive messages](#)

Name DemoQueue	Type Standard	ARN arn:aws:sqs:eu-central-1:387124123361:DemoQueue
Encryption Enabled	URL https://sns.eu-central-1.amazonaws.com/387124123361/DemoQueue	Dead-letter queue Disabled
More		

Send and receive messages

Send messages to and receive messages from a queue.

[Send message](#) [Info](#) [Clear content](#) [Send message](#)

Your message has been sent and is ready to be received. [View details](#) [X](#)

Message body
Enter the message to send to the queue.
hello world!

Delivery delay [Info](#)
0 Seconds ▾
Should be between 0 seconds and 15 minutes.

► Message attributes - Optional [Info](#)

[Receive messages](#) [Info](#) [Edit poll settings](#) [Stop polling](#) [Poll for messages](#)

Messages available 1 Polling duration 30 Maximum message count 10 Polling progress 0 receives/second

Messages (0) [View details](#) [Delete](#)

[Search messages](#)

[Poll for messages](#)

[Receive messages](#) [Info](#) [Edit poll settings](#) [Stop polling](#) [Poll for messages](#)

Messages available 1 Polling duration 30 Maximum message count 10 Polling progress 0 receives/second

Messages (0) [View details](#) [Delete](#)

[Search messages](#)

No messages. To view messages in the queue, poll for messages. [Poll for messages](#)

[Receive messages](#) [Info](#) [Edit poll settings](#) [Stop polling](#) [Poll for messages](#)

Messages available 1 Polling duration 30 Maximum message count 10 Polling progress 1 receives/second 23%

Messages (1) [View details](#) [Delete](#)

[Search messages](#)

ID	Sent	Size	Receive count
1cbe8079-fdb2-4950-a380-169c5a5f7244	02/08/2020, 16:45:45	12 bytes	1

Message: 1cbe8079-fdb2-4950-a380-169c5a5f7244

[Details](#) [Body](#) [Attributes](#)

ID 1cbe8079-fdb2-4950-a380-169c5a5f7244	Size 12 bytes	MD5 of message body fc3ff98e8c6a0d3087d51 5c0473f8677	Sender account ID 387124123361
Sent 02/08/2020, 16:45:45	First received 02/08/2020, 16:45:56	Receive count 1	Message attributes count -
Message attributes size -	MD5 of message attributes -		

[Done](#)

Delete Messages

Are you sure you want to delete the following message? You can't undo this action.

- 1cbe8079-fdb2-4950-a380-169c5a5f7244 (12 bytes)

Cancel **Delete**

Amazon SQS > Queues > DemoQueue

DemoQueue

Details **Info**

Name DemoQueue	Type Standard	ARN arn:aws:sqs:eu-central-1:387124123361:DemoQueue
Encryption Enabled	URL https://sqs.eu-central-1.amazonaws.com/387124123361/DemoQueue	Dead-letter queue Disabled
More		

Purge queue

Are you sure you want to purge the following queue permanently? **You can't undo this action.**

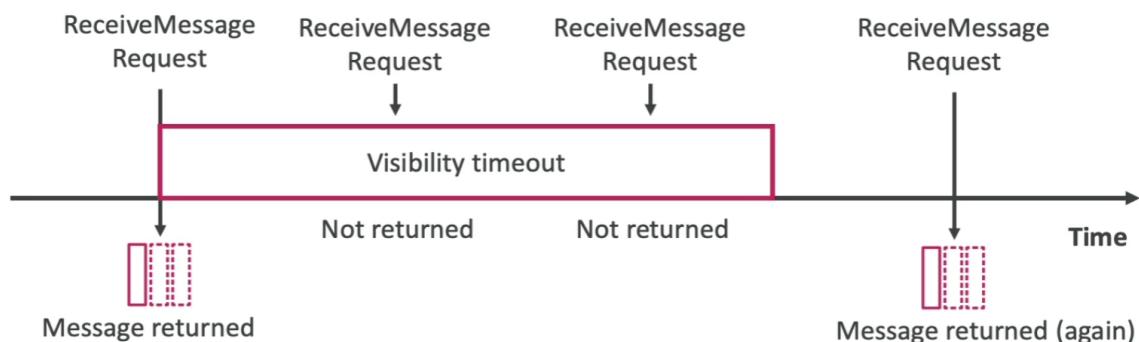
- DemoQueue - contains 0 messages

To confirm, enter the phrase **purge**.

Cancel **Purge**

SQS - Message Visibility Timeout

- After a message is polled by a consumer, it becomes invisible to other consumers
- By default, the "message visibility timeout" is 30 seconds
- That means the message has 30 seconds to be processed
- After the message visibility timeout is over, the message is "visible" in SQS

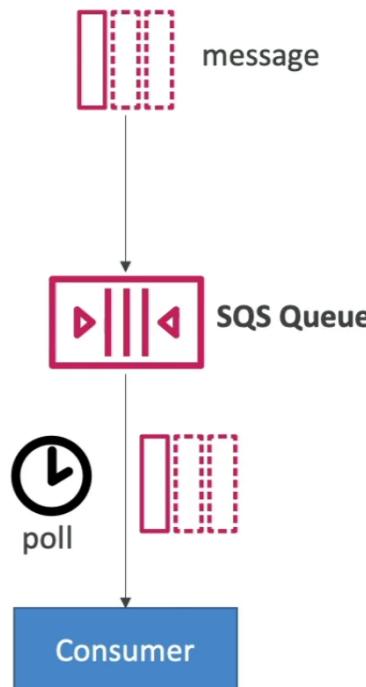


- If a message is not processed within the visibility timeout, it will be processed **twice**
- A consumer could call the **ChangeMessageVisibility** API to get more time
- If visibility timeout is high (hours), and consumer crashes, re-processing will take time
- If visibility timeout is too low (seconds), we may get duplicates

Amazon SQS - Long Polling

- When a consumer requests messages from the queue, it can optionally "wait" for messages to arrive if there are none in the queue
- This is called Long Polling
- LongPolling decreases the number of API calls made to SQS while increasing the efficiency and latency of your application
- The wait time can be between 1 sec to 20 sec (20 sec preferable)
- Long Polling is preferable to Short Polling

- Long polling can be enabled at the queue level or at the API level using **WaitTimeSeconds**



Amazon SQS - FIFO Queue

- FIFO = First In First Out (ordering of message in the queue)



- Limited throughput: 300 msg/s without batching, 3000 msg/s with
- Exactly-once send capability (by removing duplicates)
- Messages are processed in order by the consumer

Hands On

Create queue

Details

Type
Choose the queue type for your application or cloud infrastructure.

Standard Info
At-least-once delivery, message ordering isn't preserved

- At-least once delivery
- Best-effort ordering

FIFO Info
First-in-first-out delivery, message ordering is preserved

- First-in-first-out delivery
- Exactly-once processing

Name

A queue name is case-sensitive and can have up to 80 characters. You can use alphanumeric characters, hyphens (-), and underscores (_).

Configuration

Set the maximum message size, visibility to other consumers, and message retention. [Info](#)

Visibility timeout <small>Info</small>	Message retention period <small>Info</small>
<input type="text" value="30"/> Seconds	<input type="text" value="4"/> Days
Should be between 0 seconds and 12 hours.	
Delivery delay <small>Info</small>	Maximum message size <small>Info</small>
<input type="text" value="0"/> Seconds	<input type="text" value="256"/> KB
Should be between 0 seconds and 15 minutes.	
Receive message wait time <small>Info</small>	
<input type="text" value="0"/> Seconds	
Should be between 0 and 20 seconds.	

Content-based deduplication
When content-based deduplication is enabled, the message deduplication ID is optional.

Choose method

Basic
Use simple criteria to define a basic access policy.

Advanced
Use a JSON object to define an advanced access policy.

Define who can send messages to the queue

Only the queue owner
Only the owner of the queue can send messages to the queue.

Only the specified AWS accounts, IAM users and roles
Only the specified AWS account IDs, IAM users and roles can send messages to the queue.

Define who can receive messages from the queue

Only the queue owner
Only the owner of the queue can receive messages from the queue.

Only the specified AWS accounts, IAM users and roles
Only the specified AWS account IDs, IAM users and roles can receive messages from the queue.

JSON (read-only)

```
{
  "Version": "2008-10-17",
  "Id": "__default_policy_ID",
  "Statement": [
    {
      "Sid": "__owner_statement",
      "Effect": "Allow",
      "Principal": {
        "AWS": "387124123361"
      },
      "Action": [
        "SQS:*"
      ],
      "Resource": "arn:aws:sqs:eu-central-1:387124123361:DemoQueue.fifo"
    }
  ]
}
```

Receive messages Info

Messages available	Polling duration	Maximum message count	Polling progress
4	30	10	0% 0 receives/second

Messages (0)

ID	Sent	Size	Receive count
No messages. To view messages in the queue, poll for messages.			

Poll for messages

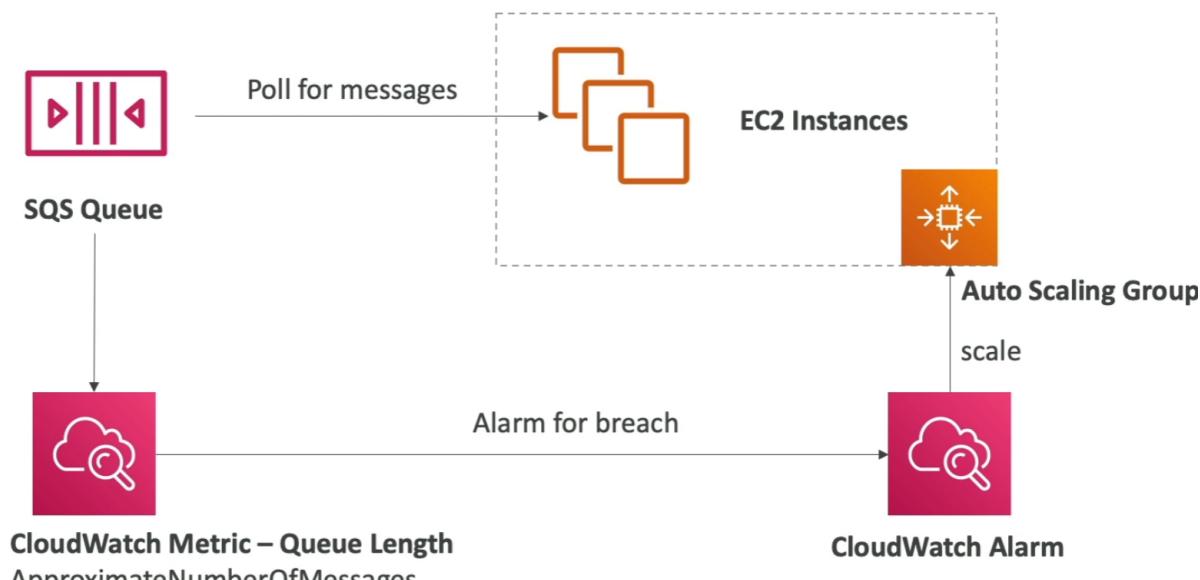
Receive messages Info

Messages available	Polling duration	Maximum message count	Polling progress
4	30	10	70% 4 receives/second

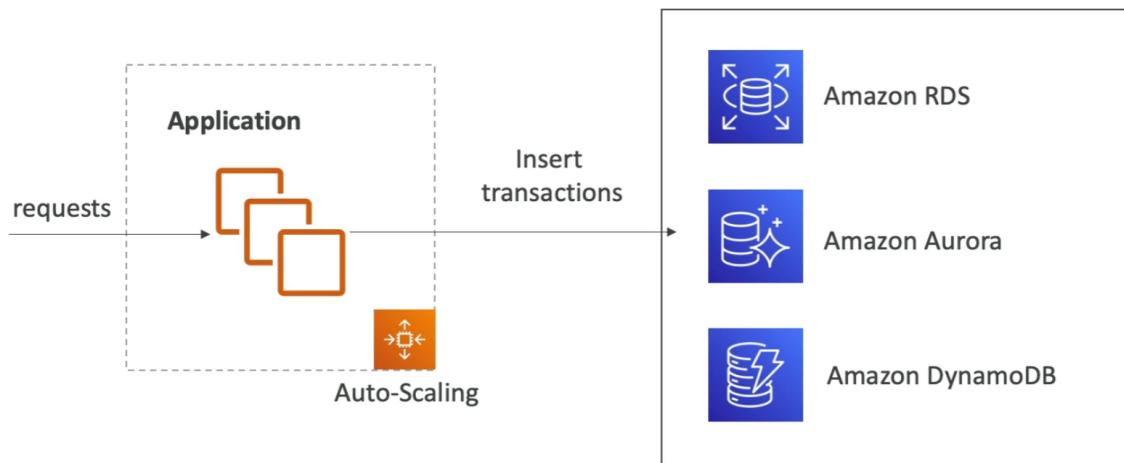
Messages (4)

ID	Sent	Size	Receive count
4768bede-0712-46d3-8c00-8d88fd3a8789	02/08/2020, 17:26:11	13 bytes	1
f0b18b8d-650-4857-899c-9b0f58bb8ff0	02/08/2020, 17:26:07	13 bytes	1
e9a3f412-c710-4ebc-af02-464dee03c91d	02/08/2020, 17:26:01	13 bytes	1
50dece5d-8870-4ad4-b30a-efb0121bca84	02/08/2020, 17:25:54	13 bytes	1

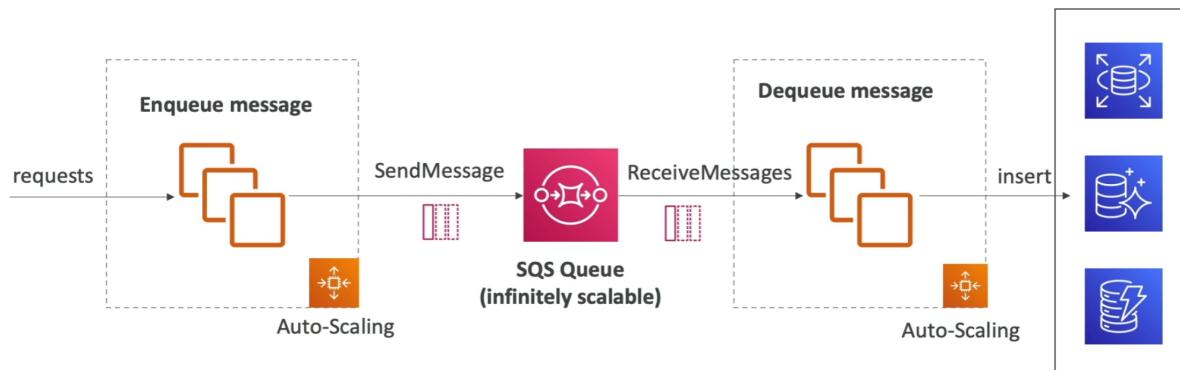
SQS with Auto Scaling Group (ASG)



If the load is too big, some transactions may be lost



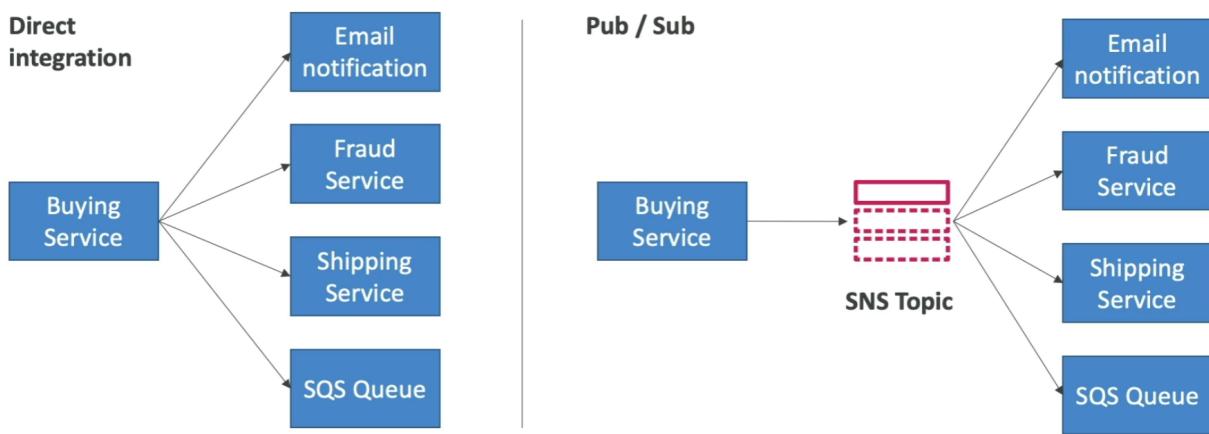
SQS as a buffer to database writes



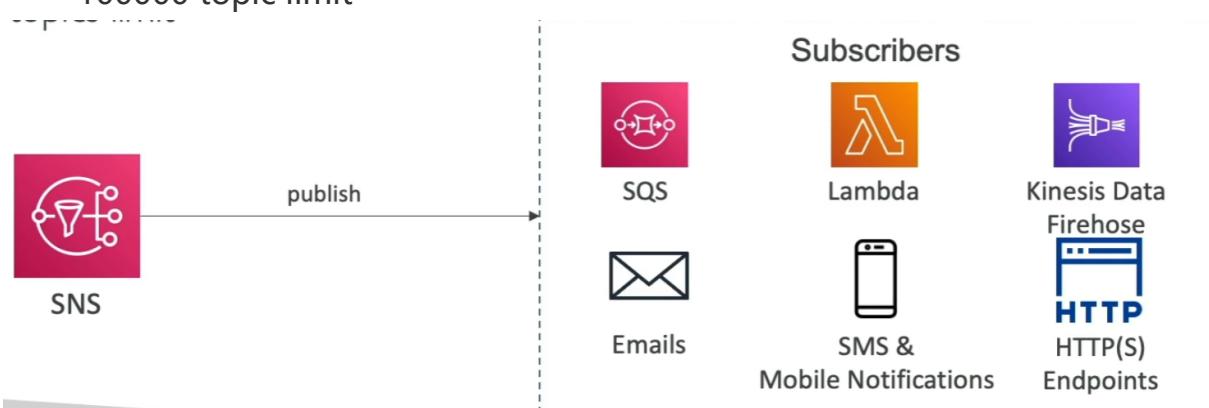
SQS to decouple between application

Amazon SNS

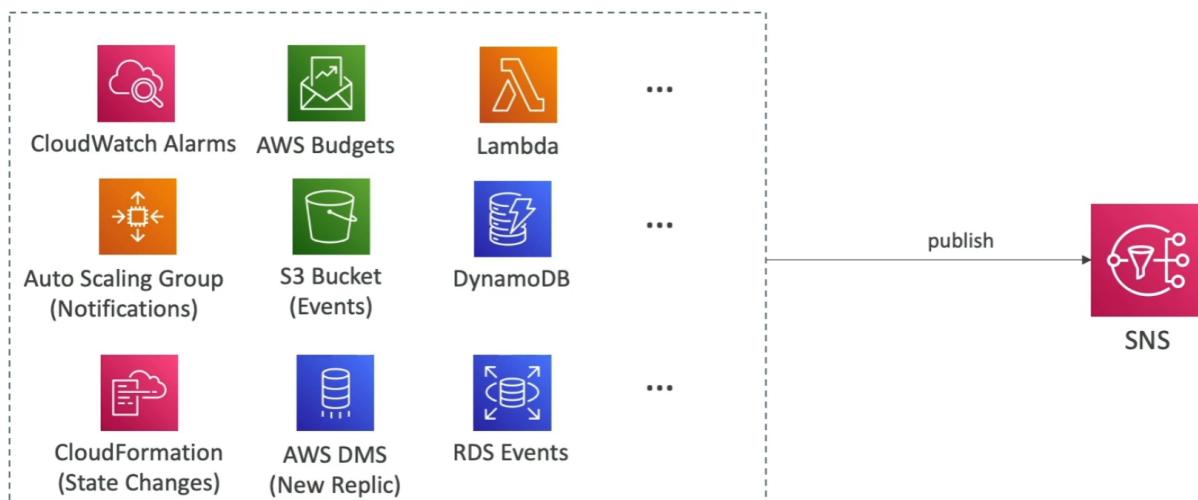
What if you want to send one message to many receivers?



- The "event producer" only sends message to one SNS topic
- As many "event receivers" (subscriptions) as we want to listen to the SNS topic notifications
- Each subscriber to the topic will get all the message (note: new feature to filter messages)
- Up to 12500000 subscriptions per topic
- 100000 topic limit



SNS integrates with a lot of AWS services



Many AWS services can send data directly to SNS for notifications

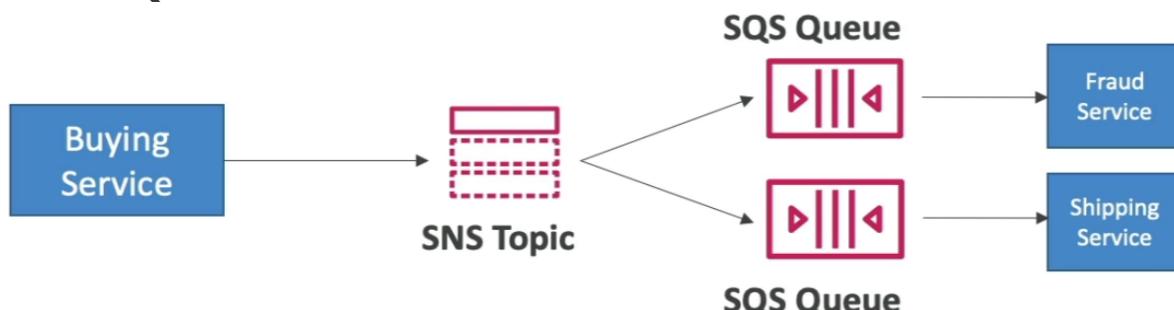
AWS SNS - How to publish

- Topic Publish (using the SDK)
 - Create a topic
 - Create a subscription (or many)
 - Publish to the topic
- Direct Publish (for mobile apps SDK)
 - Create a platform application
 - Create a platform endpoint
 - Publish to the platform endpoint
 - Works with Google GCM, Apple APNS, Amazon ADM...

Amazon SNS - Security

- **Encryption:**
 - In-flight encryption using HTTPS API
 - At-rest encryption using KMS keys
 - Client-side encryption if the client wants to perform encryption/decryption itself
- **Access Controls: IAM policies to regulate access to the SNS API**
- **SNS Access Policies (similar to S3 bucket policies)**
 - Useful for cross-account access to SNS topics
 - Useful for allowing other services (S3 ...) to write to an SNS topic

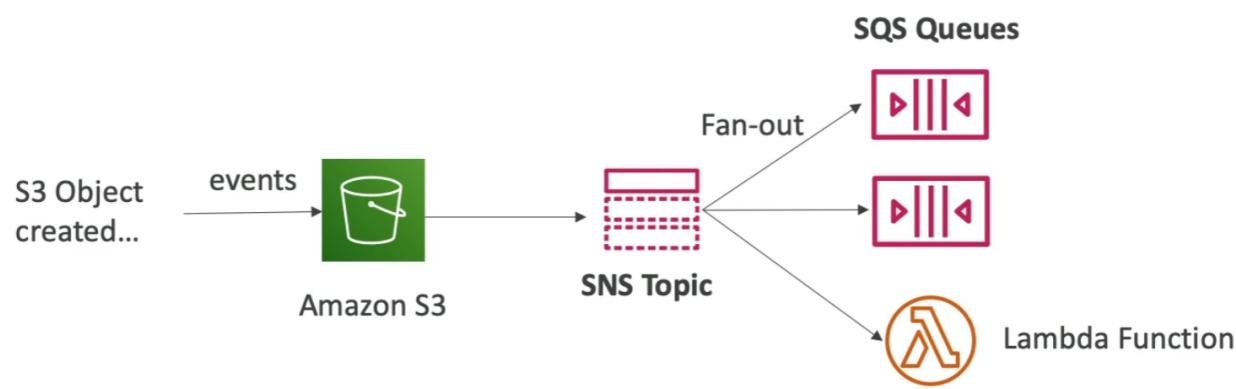
SNS + SQS: Fan Out



- Push once in SNS, receive in all SQS queues that are subscribers
- Fully decoupled, no data loss
- SQS allows for: data persistence, delayed processing and retries of work
- Ability to add more SQS subscribers over time
- Make sure your SQS queue access policy allows for SNS to write

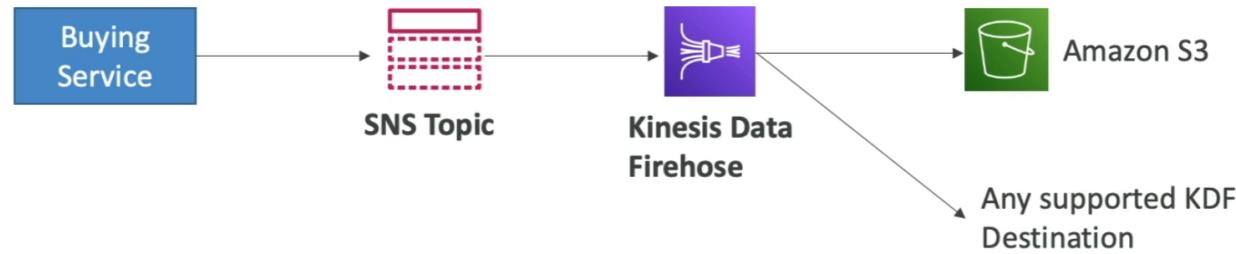
Application: S3 Events to multiple queues

- For the same combination of : event type (e.g. object create) and **prefix** (e.g. images/) you can only have one S3 Event rule
- If you want to send the same S3 event to many SQS queues, use fan-out



Application: SNS to Amazon S3 through Kinesis Data Firehose

- SNS can send to Kinesis and therefore we can have the following solutions architecture:



Amazon SNS - FIFO Topic

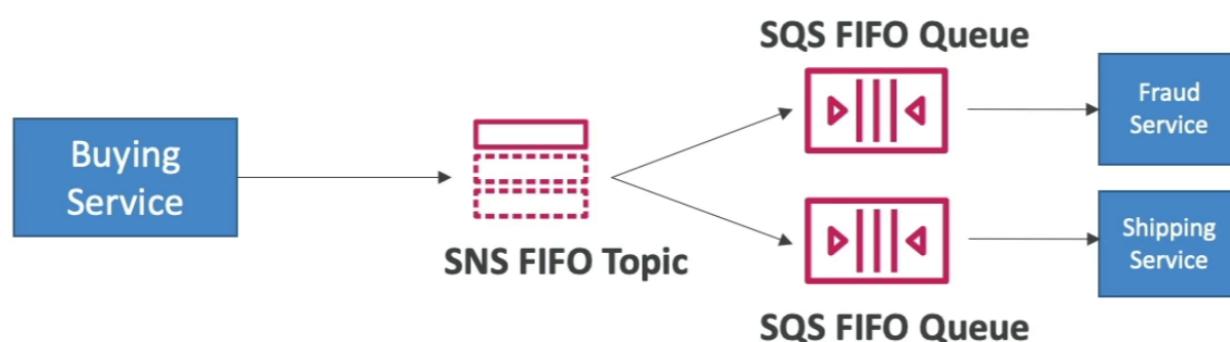
- FIFO = First In First Out (ordering of message in the topic)



- Similar features as SQS FIFO:
 - Ordering by Message Group ID (all messages in the same group are ordered)
 - Deduplication using a Deduplication ID or Content Based Deduplication
- Can only have SQS FIFO queues as subscribers
- Limited throughput (same throughput as SQS FIFO)

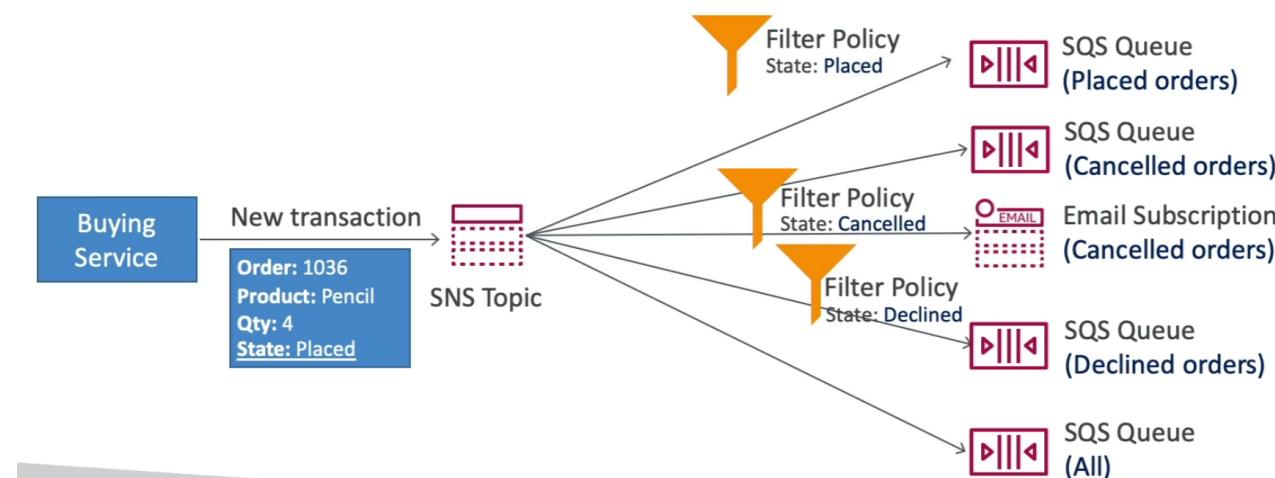
SNS FIFO + SQS FIFO: Fan Out

- In case you need fan out + ordering + deduplication



SNS - Message Filtering

- JSON policy used to filter message sent to SNS topic's subscriptions
- If a subscription doesn't have a filter policy, it receives every message



Hands On

AWS Services ▾ Search for services, features, marketplace products, and docs [Option+S] stephane-aws-course ▾ Ohio ▾

Application Integration

Amazon Simple Notification Service

Pub/sub messaging for microservices and serverless applications.

Amazon SNS is a highly available, durable, secure, fully managed pub/sub messaging service that enables you to decouple microservices, distributed systems, and event-driven serverless applications. Amazon SNS provides topics for high-throughput, push-based, many-to-many messaging.

Create topic

Topic name
A topic is a message channel. When you publish a message to a topic, it fans out the message to all subscribed endpoints.

Next step

[Start with an overview](#)

Create topic

Details

Type [Info](#)
Topic type cannot be modified after topic is created

FIFO (first-in, first-out)

- Strictly-preserved message ordering
- Exactly-once message delivery
- High throughput, up to 300 publishes/second
- Subscription protocols: SQS

Standard

- Best-effort message ordering
- At-least once message delivery
- Highest throughput in publishes/second
- Subscription protocols: SQS, Lambda, HTTP, SMS, email, mobile application endpoints

Name

Maximum 256 characters. Can include alphanumeric characters, hyphens (-) and underscores (_).

Display name - *optional*
To use this topic with SMS subscriptions, enter a display name. Only the first 10 characters are displayed in an SMS message. [Info](#)

Maximum 100 characters, including hyphens (-) and underscores (_).

► Encryption - *optional*
Amazon SNS provides in-transit encryption by default. Enabling server-side encryption adds at-rest encryption to your topic.

► Access policy - *optional*
This policy defines who can access your topic. By default, only the topic owner can publish or subscribe to the topic. [Info](#)

► Delivery retry policy (HTTP/S) - *optional*
The policy defines how Amazon SNS retries failed deliveries to HTTP/S endpoints. To modify the default settings, expand this section. [Info](#)

► Delivery status logging - *optional*
These settings configure the logging of message delivery status to CloudWatch Logs. [Info](#)

► Tags - *optional*
A tag is a metadata label that you can assign to an Amazon SNS topic. Each tag consists of a key and an optional value. You can use tags to search and filter your topics and track your costs. [Learn more](#)

[Cancel](#) **Create topic**

Subscriptions (0)					
Edit	Delete	Request confirmation	Confirm subscription	Create subscription	
<input type="text" value="Search"/>					
ID	Endpoint	Status	Protocol		
No subscriptions found					
You don't have any subscriptions to this topic.					
Create subscription					

Create subscription

Details

Topic ARN
 X

Protocol
The type of endpoint to subscribe
 ▼

Endpoint
An email address that can receive notifications from Amazon SNS.

ⓘ After your subscription is created, you must confirm it. [Info](#)

► **Subscription filter policy - optional**
This policy filters the messages that a subscriber receives. [Info](#)

► **Redrive policy (dead-letter queue) - optional**
Send undeliverable messages to a dead-letter queue. [Info](#)

[Cancel](#) Create subscription

MyFirstTopic

Details

Name MyFirstTopic	Display name -
ARN arn:aws:sns:us-east-2:001736599714:MyFirstTopic	Topic owner 001736599714
Type Standard	

Message body

Message structure

Identical payload for all delivery protocols.
The same payload is sent to endpoints subscribed to the topic, regardless of their delivery protocol.

Custom payload for each delivery protocol.
Different payloads are sent to endpoints subscribed to the topic, based on their delivery protocol.

Message body to send to the endpoint

```
1 hello world
```

To stephanetheteacher	Delete
From no-reply@sns.amazonaws.com	
Sending 54.240.30.9	
IP	
Received 2021-04-05 23:13:06	

TEXT **JSON** **RAW** **LINKS** **ATTACHMENTS**

hello world

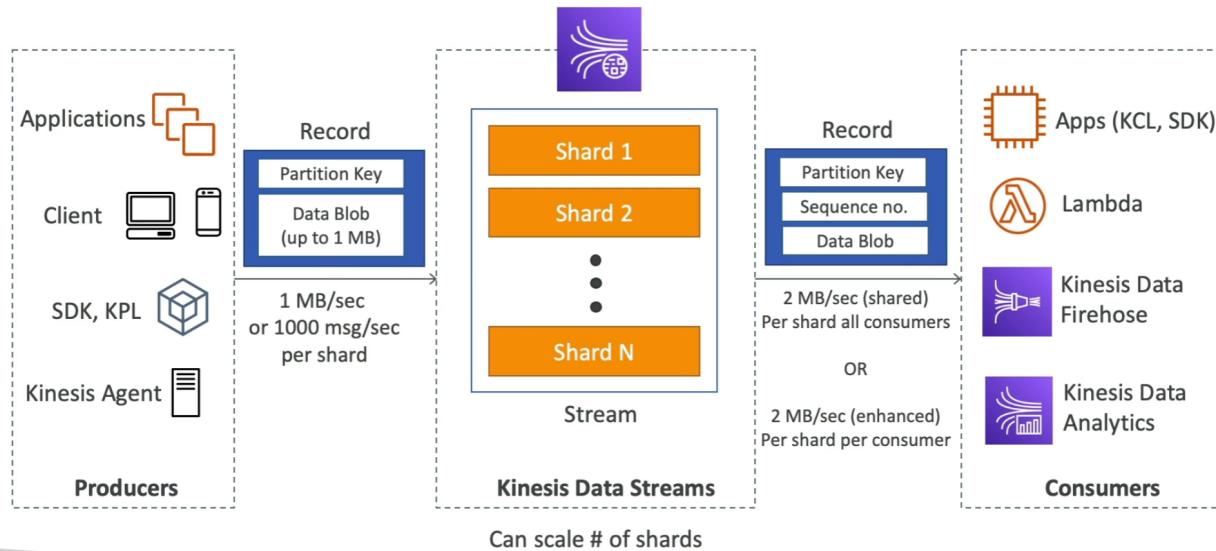
--

If you wish to stop receiving notifications from this topic, please click or visit the link below to unsubscribe:
<https://sns.us-east-2.amazonaws.com/unsubscribe.html?SubscriptionArn=arn:aws:sns:us-east-2:001736599714:MyFirstTopic:57aa54d4-e48b-4d55-a28f-bc3292693323&Endpoint=stephanetheteacher@mailinator.com>

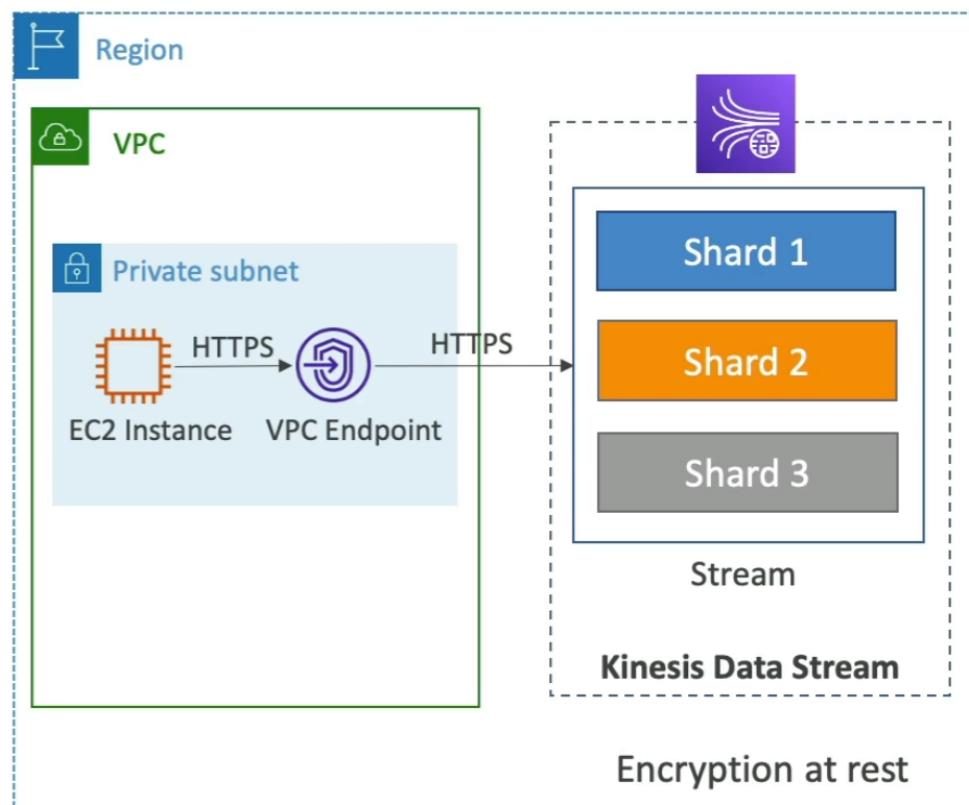
Please do not reply directly to this email. If you have any questions or comments regarding this email, please contact us at
<https://aws.amazon.com/support>

Kinesis Overview

- Makes it easy to collect, process, and analyze streaming data in real-time
- Ingest real-time data such as: Application logs, Metrics, Website clickstreams, IoT telemetry data ...
- **Kinesis Data Streams:** capture, process, and store data streams



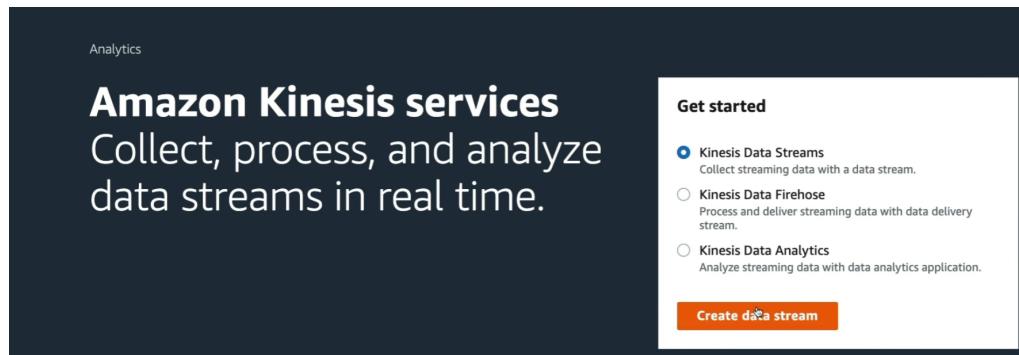
- Retention between 1 day to 365 days
- Ability to reprocess (replay) data
- Once data is inserted in Kinesis, it can't be deleted (immutability)
- Data that shares the same partition goes to the same shard (ordering)
- Producers: AWS SDK, Kinesis Producer Library (KPL), Kinesis Agent
- Consumers:
 - write your own: Kinesis Client Library(KCL),AWS SDK
 - Managed:AWS Lambda, Kinesis Data Firehose, Kinesis Data Analytics
- Capacity Modes
 - **Provisioned mode:**
 - You choose the number of shards provisioned, scale manually or using API
 - Each shard gets 1MB/s in (or 1000 records per second)
 - Each shard gets 2MB/s out (classic or enhanced fan-out consumer)
 - You pay per shard provisioned per hour
 - **On-demand mode:**
 - No need to provision or manage the capacity
 - Default capacity provisioned (4 MB/s in or 4000 records per second)
 - Scales automatically based on observed throughput peak during the last 30 days
 - Pay per stream per hour & data in/out per GB
- Security



- Control access / authorization using IAM policies
- Encryption in flight using HTTPS endpoints
- Encryption at rest using KMS
- You can implement encryption/decryption of data on client side (harder)
- VPC Endpoints available for Kinesis to access within VPC

- Monitor API calls using CloudTrail

- Hands on



Create data stream Info

Data stream configuration

Data stream name
DemoStream
Acceptable characters are uppercase and lowercase letters, numbers, underscores, hyphens and periods.

Data stream capacity Info

Capacity mode

- On-demand**
Use this mode when your data stream's throughput requirements are unpredictable and variable. With on-demand mode, your data stream's capacity scales automatically.
- Provisioned**
Use provisioned mode when you can reliably estimate throughput requirements of your data stream. With provisioned mode, your data stream's capacity is fixed.

Provisioned shards
The total capacity of a stream is the sum of the capacities of its shards. Enter number of provisioned shards to see total data stream capacity.
1
Minimum: 1, Maximum available: 200, Account quota limit: 200. [Request shard quota increase](#)

Total data stream capacity
Shard capacity is determined by the number of provisioned shards. Each shard ingests up to 1 MiB/second and 1,000 records/second and emits up to 2 MiB/second. If writes and reads exceed capacity, the application will receive throttles.

Write capacity Maximum 1 MiB/second and 1,000 records/second	Read capacity Maximum 2 MiB/second
---	---

Provisioned mode has a fixed-throughput pricing model. See [Kinesis pricing for Provisioned mode](#)

Data stream settings
You can edit the settings after the data stream has been created and is in the active status.

Setting	Value	Editable after creation
Capacity mode	Provisioned	<input checked="" type="checkbox"/> Yes
Provisioned shards	1	<input checked="" type="checkbox"/> Yes
Data retention period	1 day	<input checked="" type="checkbox"/> Yes
Server-side encryption	Disabled	<input checked="" type="checkbox"/> Yes
Monitoring enhanced metrics	Disabled	<input checked="" type="checkbox"/> Yes
Tags	-	<input checked="" type="checkbox"/> Yes

Create data stream

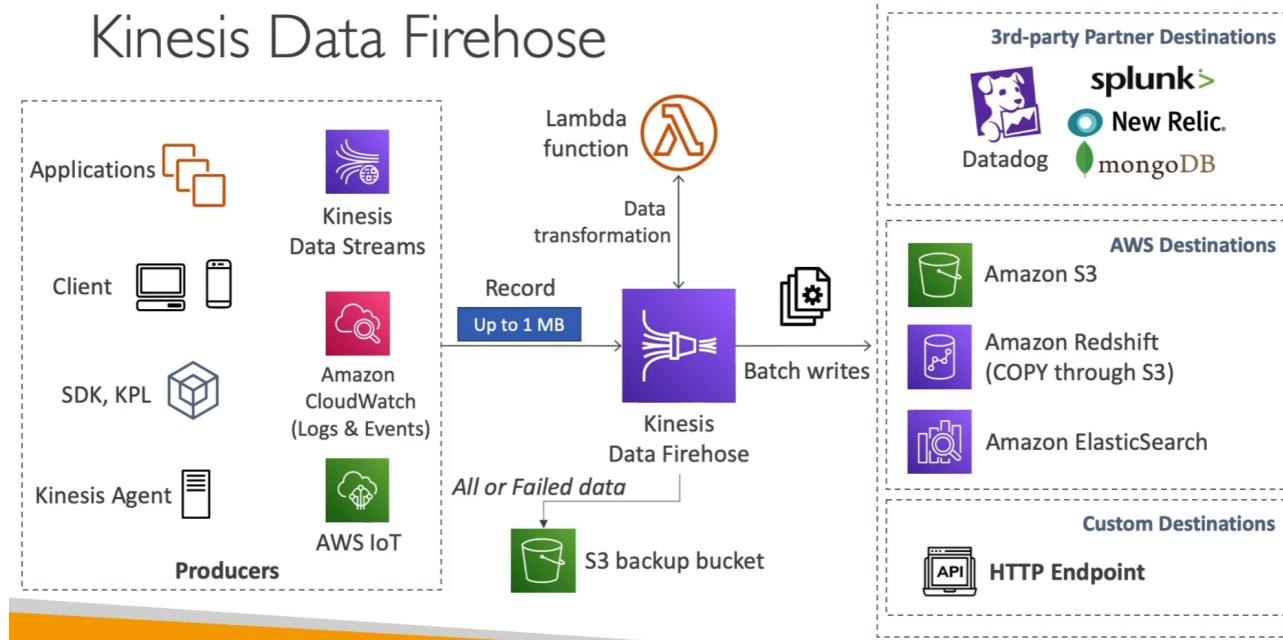
```
[cloudshell-user@ip-10-0-39-58 ~]$ Try these commands to get started:  
aws help or aws <commands> help or aws <command> --cli-auto-prompt  
[cloudshell-user@ip-10-0-39-58 ~]$ aws --version  
aws-cli/2.1.16 Python/3.7.3 Linux/4.14.225-168.357.amzn2.x86_64 exec-env/CloudShell.exe/x86_64.amzn.2 prompt/off  
[cloudshell-user@ip-10-0-39-58 ~]$ aws kinesis put-record --stream-name DemoStream  
--partition-key user1 --data "user signup" --cli-binary-format raw-in-base64-out  
{  
    "ShardId": "shardId-000000000000",  
    "SequenceNumber": "49617390934629201455926329624637508189499828916235272194"  
}  
[cloudshell-user@ip-10-0-39-58 ~]$ aws kinesis put-record --stream-name DemoStream  
--partition-key user1 --data "user signup" --cli-binary-format raw-in-base64-out  
{  
    "ShardId": "shardId-000000000000",  
    "SequenceNumber": "49617390934629201455926329624638717115319444438763175938"  
}  
[cloudshell-user@ip-10-0-39-58 ~]$ aws kinesis put-record --stream-name DemoStream  
--partition-key user1 --data "user login" --cli-binary-format raw-in-base64-out  
{  
    "ShardId": "shardId-000000000000",  
    "SequenceNumber": "49617390934629201455926329624639926041139059686413172738"  
}  
[cloudshell-user@ip-10-0-39-58 ~]$ aws kinesis put-record --stream-name DemoStream  
--partition-key user1 --data "user logout" --cli-binary-format raw-in-base64-out
```

```
[cloudshell-user@ip-10-0-39-58 ~]$ aws kinesis describe-stream --stream-name DemoStream  
{  
    "StreamDescription": {  
        "Shards": [  
            {  
                "ShardId": "shardId-000000000000",  
                "HashKeyRange": {  
                    "StartingHashKey": "0",  
                    "EndingHashKey": "340282366920938463463374607431768211455"  
                },  
                "SequenceNumberRange": {  
                    "StartingSequenceNumber": "49617390934629201455926329620716961756489569047916052482"  
                }  
            }  
        ],  
        "StreamARN": "arn:aws:kinesis:us-east-1:001736599714:stream/DemoStream",  
        "StreamName": "DemoStream",  
        "StreamStatus": "ACTIVE",  
        "RetentionPeriodHours": 24,  
        "EnhancedMonitoring": [  
            {  
                "ShardLevelMetrics": []  
            }  
        ]  
    }  
}
```

```
[cloudshell-user@ip-10-0-39-58 ~]$ aws kinesis get-shard-iterator --stream-name DemoStream --shard-id shardId-000000000000 --shard-iterator-type TRIM_HORIZON  
{  
    "ShardIterator": "AAAAAAAAGtYgUq00wFKpIqqDHDoSwgADDT3VzUW6qPzgR04TPuWEbCgUZPVjI2adlKwjAUssIIpjLY0aC82nhF+ML2nLT5WxK3WLltuh340ZrsVw9Q0gT6QpaZlkqi6wSKfb2yrqxj6ShbHueiUuqQMH/5hW2NRaLhgdxm6thQtoZF+BS+ljcye7BJV1iwNKz29KH96EqGnn9UcxLyZoSLWxD/KV1uck88QV+u02S/HQ1+v/xKmw=="  
}  
[cloudshell-user@ip-10-0-39-58 ~]$ aws kinesis get-records --shard-iterator "AAAAAAAGtYgUq00wFKpIqqDHDoSwgADDT3VzUW6qPzgR04TPuWEbCgUZPVjI2adlKwjAUssIIpjLY0aC82nhF+ML2nLT5WxK3WLltuh340ZrsVw9Q0gT6QpaZlkqi6wSKfb2yrqxj6ShbHueiUuqQMH/5hW2NRaLhgdxm6thQtoZF+BS+ljcye7BJV1iwNKz29KH96EqGnn9UcxLyZoSLWxD/KV1uck88QV+u02S/HQ1+v/xKmw=="
```

```
{  
    "Records": [  
        {  
            "SequenceNumber": "49617390934629201455926329624637508189499828916235272194",  
            "ApproximateArrivalTimestamp": "2021-04-16T17:30:24.565000+00:00",  
            "Data": "dXNlcjBzaWdudXA=",  
            "PartitionKey": "user1"  
        },  
        {  
            "SequenceNumber": "49617390934629201455926329624638717115319444438763175938",  
            "ApproximateArrivalTimestamp": "2021-04-16T17:30:37.358000+00:00",  
            "Data": "dXNlcjBzaWdudXA=",  
            "PartitionKey": "user1"  
        },  
        {  
            "SequenceNumber": "49617390934629201455926329624639926041139059686413172738",  
            "ApproximateArrivalTimestamp": "2021-04-16T17:30:45.978000+00:00",  
            "Data": "dXNlcjBsb2dpbg==",  
            "PartitionKey": "user1"  
        },  
        {  
            "SequenceNumber": "49617390934629201455926329624639926041139059686413172739",  
            "ApproximateArrivalTimestamp": "2021-04-16T17:30:46.000000+00:00",  
            "Data": "dXNlcjBsb2dpbg==",  
            "PartitionKey": "user1"  
        }  
    ]
```

- **Kinesis Data Firehose:** load data streams into AWS data stores



- Fully Managed Service, no administration, automatic scaling. serverless
 - AWS: Redshift / Amazon S3 / ElasticSearch
 - 3rd party partner: Splunk / MongoDB / DataDog / NewRelic /...
 - Custom: send to any HTTP endpoint
- Pay for data going through Firehose
- **Near Real Time**
 - 60 seconds latency minimum for non full batches
 - Or minimum 1 MB of data at a time
- Supports many data formats, conversions, transformations, compression
- Supports custom data transformations using AWS Lambda
- Can send failed or all data to a backup S3 bucket

- Hands on

Amazon Kinesis

Amazon Kinesis > Data streams > DemoStream

DemoStream [Info](#)

[Delete](#)

Data stream summary			
Status Active	Data retention period 1 day	ARN arn:aws:kinesis:eu-central-1:211442049068:stream/DemoStream	Creation time October 21, 2021, 13:25 GMT+1

[Applications](#) [Monitoring](#) [Configuration](#) [Enhanced fan-out \(0\)](#)

Amazon Kinesis

Amazon Kinesis > Delivery streams

Delivery streams

[Create delivery stream](#)

Name	Status	Creation time	Source	Data transformation	Destination	Delivery status
No delivery streams No delivery streams to display						

[Create delivery stream](#)

Choose source and destination

Specify the source and the destination for your delivery stream. You cannot change the source and destination of your delivery stream once it has been created.

Source [Info](#)
Amazon Kinesis Data Streams

Destination [Info](#)
Amazon S3

Source settings

Kinesis data stream
arn:aws:kinesis:eu-central-1:211442049068:stream/DemoStream [Browse](#) [Create](#)

Format: arn:aws:kinesis:[Region]:[AccountId]:stream/[StreamName]

Delivery stream name

Delivery stream name
KDS-S3-BRnsY

Acceptable characters are uppercase and lowercase letters, numbers, underscores, hyphens, and periods.

Transform and convert records - optional

Configure Kinesis Data Firehose to transform and convert your record data.

Transform source records with AWS Lambda [Info](#)
Kinesis Data Firehose can invoke an AWS Lambda function to transform, filter, un-compress, convert and process your source data records. The specified AWS Lambda function can also be used to provide dynamic partitioning keys for the incoming source data before its delivery to the specified destination.

Data transformation
 Disabled
 Enabled

Convert record format [Info](#)
Data in Apache Parquet or Apache ORC format is typically more efficient to query than JSON. Kinesis Data Firehose can convert your JSON-formatted source records using a schema from a table defined in [AWS Glue](#). For records that aren't in JSON format, create a Lambda function that converts them to JSON in the Transform source records with AWS Lambda section above.

Record format conversion
 Disabled
 Enabled

Destination settings [Info](#)

Specify the destination settings for your delivery stream.

S3 bucket

[Browse](#)
[Create](#)

Format: s3://bucket

Dynamic partitioning [Info](#)

Dynamic partitioning enables you to create targeted data sets by partitioning streaming S3 data based on partitioning keys. You can partition your source data with inline parsing and/or the specified AWS Lambda function. You can enable dynamic partitioning only when you create a new delivery stream. You cannot enable dynamic partitioning for an existing delivery stream. Enabling dynamic partitioning incurs additional costs per GiB of partitioned data. For more information, see [Kinesis Data Firehose pricing](#).

Disabled

Enabled

S3 bucket prefix - optional

By default, Kinesis Data Firehose appends the prefix "YYYY/MM/dd/HH" (in UTC) to the data it delivers to Amazon S3. You can override this default by specifying a custom prefix that includes expressions that are evaluated at runtime.

Enter a prefix

Amazon CloudWatch error logging [Info](#)

Choose Enabled if you want Kinesis Data Firehose to log record delivery errors to CloudWatch Logs.

Disabled

Enabled

Permissions [Info](#)

Kinesis Data Firehose uses this IAM role for all the permissions that the delivery stream needs. To specify different roles for the different permissions, use the API or the CLI.

Create or update IAM role KinesisFirehoseServiceRole-KDS-S3-BRn-eu-central-1-1634819284161

Creates a new role or updates an existing one and adds the required policies to it, and enables Kinesis Data Firehose to assume it.

Choose existing IAM role

The role that you choose must have policies that include the permissions that Kinesis Data Firehose needs.

Tags [Info](#)

You can add tags to organize your AWS resources, track costs, and control access.

No tags associated with the resource.

[Add new tag](#)

You can add up to 50 more tags.

[Cancel](#)
[Create delivery stream](#)

```
[cloudshell-user@ip-10-0-38-107 ~]$ aws kinesis put-record --stream-name DemoStream --partition-key user1 --data "user signup" --cli-binary-format raw-in-base64-out
{
  "ShardId": "shardId-000000000000",
  "SequenceNumber": "49623180259820468866284917982806735146293046538896670722"
}
[cloudshell-user@ip-10-0-38-107 ~]$ aws kinesis put-record --stream-name DemoStream --partition-key user1 --data "user login" --cli-binary-format raw-in-base64-out
{
  "ShardId": "shardId-000000000000",
  "SequenceNumber": "49623180259820468866284917982807944072112661580388237314"
}
[cloudshell-user@ip-10-0-38-107 ~]$ aws kinesis put-record --stream-name DemoStream --partition-key user1 --data "user logout" --cli-binary-format raw-in-base64-out
{
  "ShardId": "shardId-000000000000",
  "SequenceNumber": "49623180259820468866284917982809152997932276553160327170"
}
[cloudshell-user@ip-10-0-38-107 ~]$
```

Amazon S3 > demo-firehose-stephane-v3

demo-firehose-stephane-v3 [Info](#)

[Objects](#) [Properties](#) [Permissions](#) [Metrics](#) [Management](#) [Access Points](#)

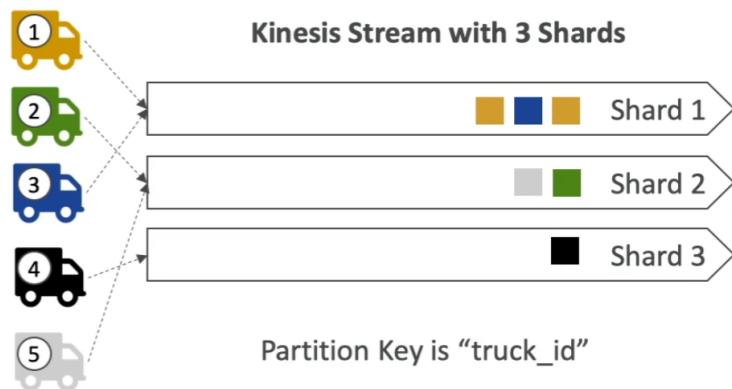
Objects (1)																				
Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more																				
 Actions 																				
 																				
<input type="text"/> Find objects by prefix																				
<table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="width: 10px;"></th> <th>Name</th> <th>Type</th> <th>Last modified</th> <th>Size</th> <th>Storage class</th> <th style="width: 10px;"></th> </tr> </thead> <tbody> <tr> <td><input type="checkbox"/></td> <td>2021/</td> <td>Folder</td> <td>-</td> <td>-</td> <td>-</td> <td style="text-align: right;">-</td> </tr> </tbody> </table>								Name	Type	Last modified	Size	Storage class		<input type="checkbox"/>	2021/	Folder	-	-	-	-
	Name	Type	Last modified	Size	Storage class															
<input type="checkbox"/>	2021/	Folder	-	-	-	-														

Kinesis Data Streams vs Firehose

Kinesis Data Streams	Kinesis Data Firehose
<ul style="list-style-type: none"> Streaming service for ingest at scale Write custom code (producer / consumer) Real-time (~200ms) Manage scaling (shard splitting / merging) Data storage for 1 to 365 days Supports replay capability 	<ul style="list-style-type: none"> Load streaming data into S3 / Redshift / ES / 3rd party / custom HTTP Fully managed Near real-time (buffer time min. 60 sec) Automatic scaling No data storage Doesn't support replay capability

- Kinesis Data Analytics:** analyze data streams with SQL or Apache Flink
- Kinesis Video Streams:** capture, process, and store video streams

Ordering data into Kinesis



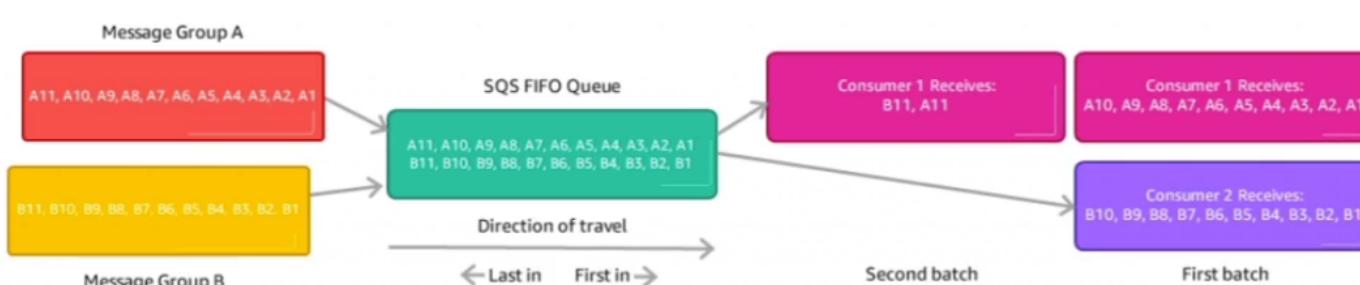
- Imagine you have 100 trucks (truck_1, truck_2, ..., truck_100) on the road sending their GPS positions regularly into AWS
- You want to consume the data in order for each truck, so that you can track their movement accurately.
- How should you send that data into Kinesis?
- Answer: send using a "Partition Key" value of the "truck_id"**
- The same key will always go to the same shard

Ordering data into SQS

- For SQS standard, there is no ordering.
- For SQS FIFO , if you don't use a Group ID, messages are consumed in the order they are sent, **with only one consumer**



- You want to scale the number of consumers, but you want messages to be "grouped" when they are related to each other
- Then you use a Group ID (similar to Partition Key in Kinesis)



Kinesis vs SQS ordering

- Let's assume 100 trucks, 5 kinesis shards, 1 SQS FIFO**
- Kinesis Data Streams:
 - On average you'll have 20 trucks per shard
 - Trucks will have their data ordered within each shard

- The maximum amount of consumers in parallel we can have is 5
- Can receive up to 5 MB/s of data
- SQS FIFO
 - You only have one SQS FIFO queue
 - You will have 100 Group ID
 - You can have up to 100 Consumers (due to the 100 Group ID)
 - You have up to 300 messages per second (or 3000 if using batching)

SQS vs SNS vs Kinesis

SQS	SNS	Kinesis
<ul style="list-style-type: none"> • Consumer "pull data" • Data is deleted after being consumed • Can have as many workers(consumers) as we want • No need to provision throughput • Ordering guarantees only on FIFO queues • Individual message delay capability 	<ul style="list-style-type: none"> • Push data to many subscribers • Up to 12500000 subscribers • Data is not persisted (lost if not delivered) • Pub/Sub • Up to 100000 topics • No need to provision throughput • Integrates with SQS for fan-out architecture pattern • FIFO capability for SQS FIFO 	<ul style="list-style-type: none"> • Standard: pull data <ul style="list-style-type: none"> ◦ 2MB per shard • Enhanced-fan out: push data <ul style="list-style-type: none"> ◦ 2 MB per shard per consumer • Possibility to replay data • Meant for real-time big data, analytics and ETL • Ordering at the shard level • Data expires after X days • Provisioned mode or on-demand capacity mode

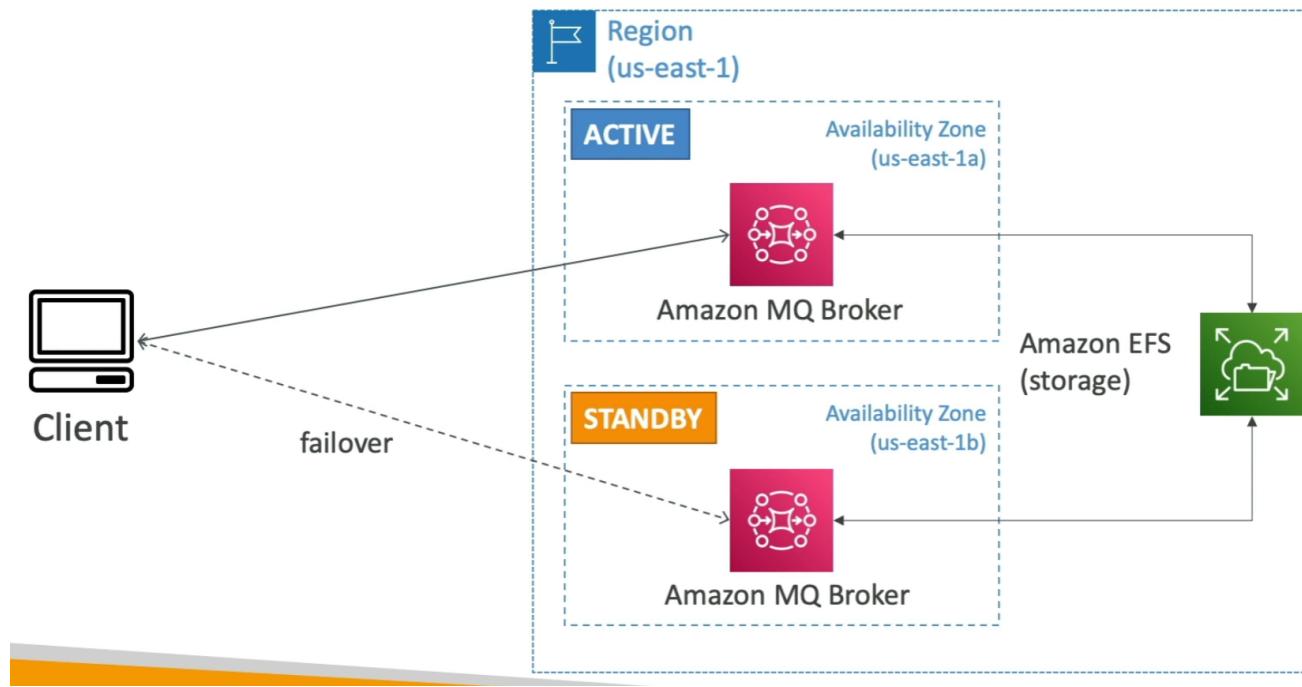
Amazon MQ

- SQS, SNS are "cloud-native" services: proprietary protocols form AWS
- Traditional applications running from on-premises may use open protocols such as: MQTT, AMQP, STOMP, Openwrite, WSS
- When migrating to the cloud, instead of re-engineering the application to use SQS and SNS, we can use Amazon MQ
- Amazon MQ is a managed message broker service for



- Amazon MQ doesn't "scale" as much as SQS / SNS
- Amazon MQ runs on servers, can run in Multi-AZ with failover
- Amazon MQ has both queue feature (~SQS) and topic features (~SNS)

Amazon MQ - High Availability



Amazon Elastic Container Service (Amazon ECS)

Amazon's own container platform

Amazon Elastic Kubernetes Service (Amazon EKS)

Amazon's managed Kubernetes(open source)

AWS Fargate

Amazon's own Serverless container platform

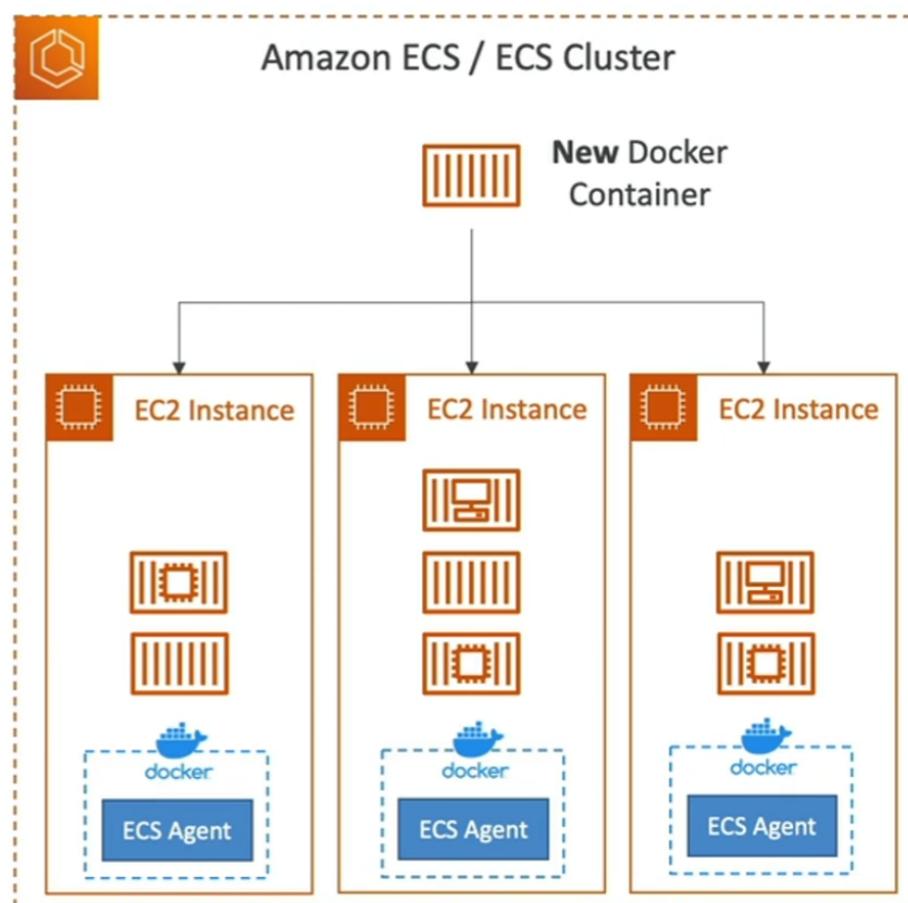
Works with ECS and with EKS

Amazon ECR:

Store container images

Amazon ECS - EC2 Launch Type

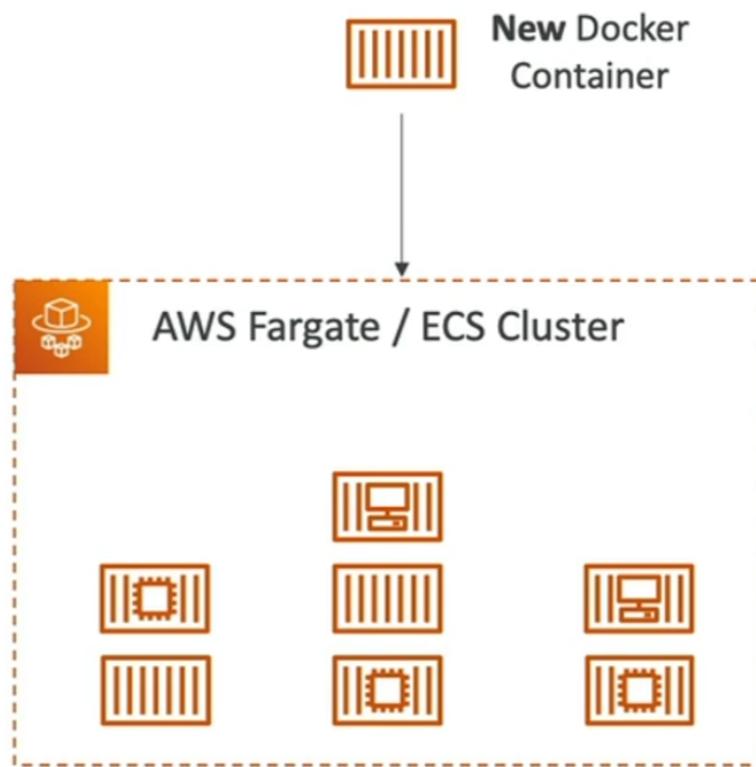
- ECS = Elastic Container Service
- Launch Docker container on AWS = Launch ECS Tasks on ECS Clusters
- EC2 Launch Type: you must provision & maintain the infrastructure (the EC2 instances)
- Each EC2 Instance must run the ECS Agent to register in the ECS Cluster
- AWS takes care of starting / stopping containers



Amazon ECS - Fargate Launch Type

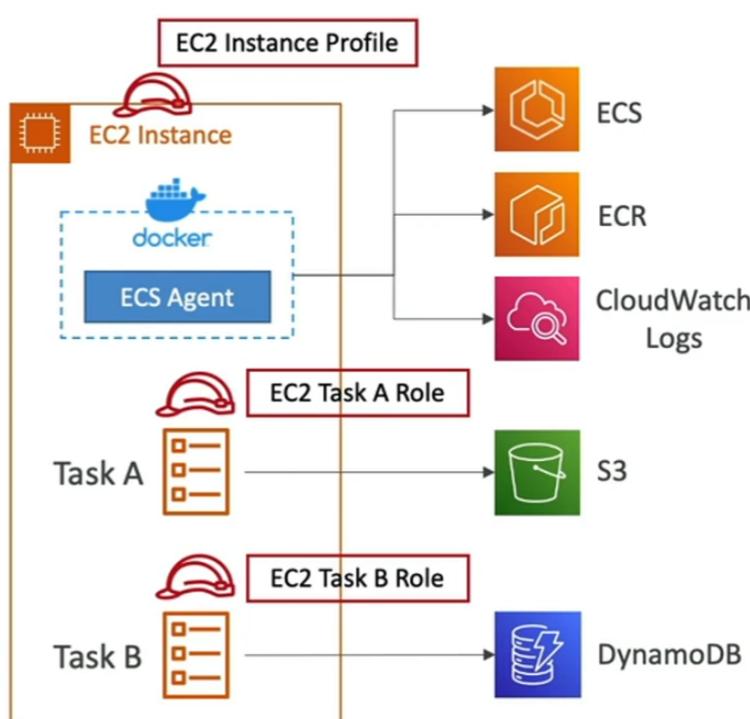
- Launch Docker containers on AWS
- You do not provision the infrastructure (no EC2 instances to manage)
- It's all Serverless!

- You just create task definitions
- AWS just runs ECS Tasks for you based on the CPU / RAM you need
- To scale, just increase the number of task. Simple - no more EC2 instances



Amazon ECS - IAM Roles for ECS

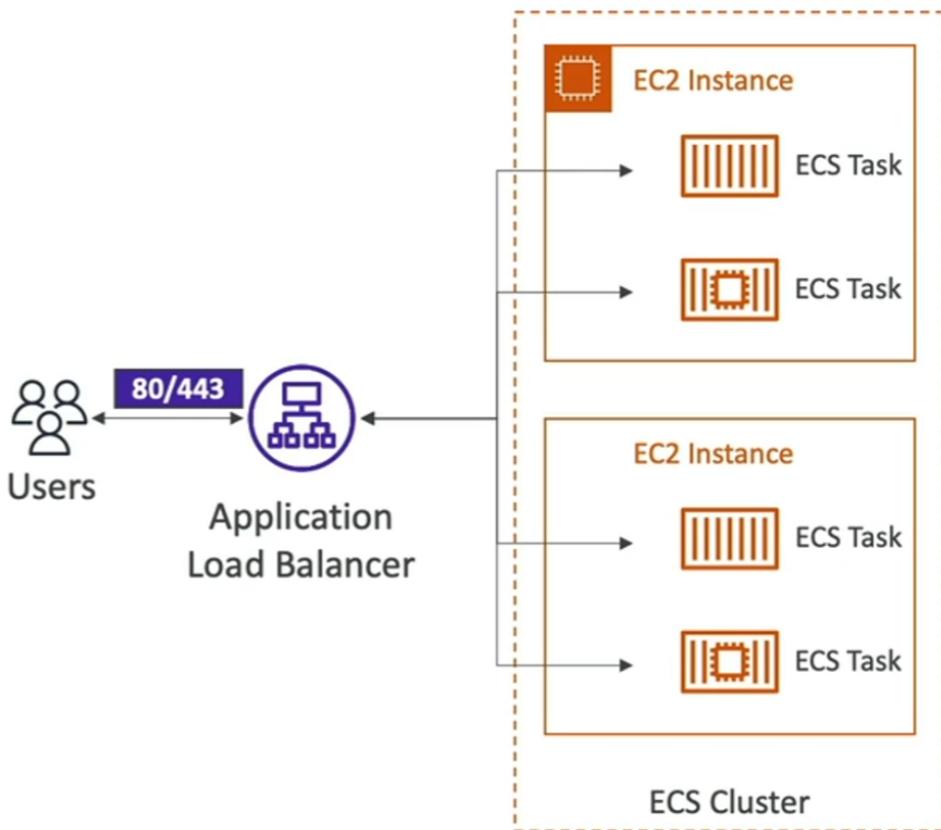
- EC2 Instance Profile (EC2 Launch Type only):
 - Used by the ECS agent
 - Makes API calls to ECS service
 - Send container logs to CloudWatch Logs
 - Pull Docker image from ECR
 - Reference sensitive data in Secrets Manager or SSM Parameter Store
- ECS Task Role:
 - Allows each task to have a specific role
 - Use different roles for the different ECS Services you run
 - Task Role is defined in the task definition
-



Amazon ECS - Load Balancer Integrations

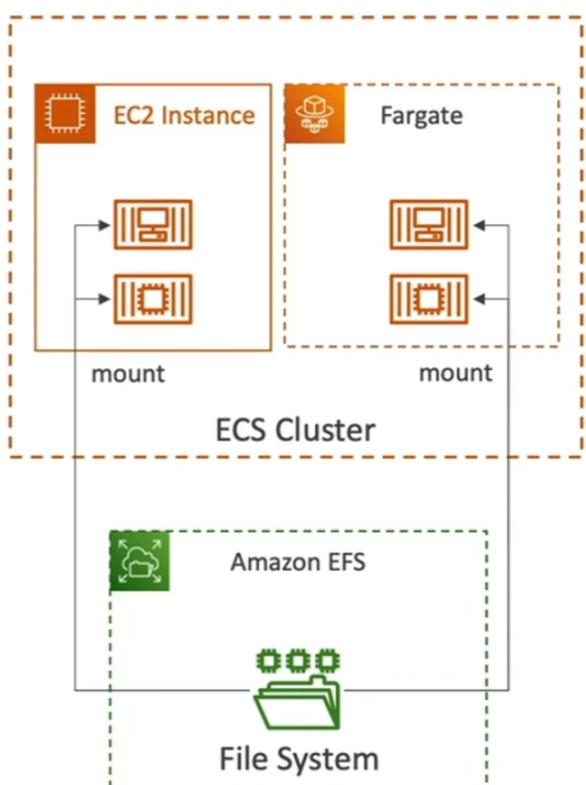
- [Application Load Balancer](#) supported and works for most use cases
- [Network Load Balancer](#) recommended only for high throughput / high performance use cases, or to pair it with AWS Private Link

- **Elastic Load Balancer** supported but not recommended (no advanced features - no Fargate)



Amazon ECS - Data Volumes (EFS)

- Mount EFS file systems onto ECS tasks
- Works for both **EC2** and **Fargate** launch types
- Tasks running in any AZ will share the same data in the EFS file system
- **Fargate + EFS = Serverless**
- Use cases: persistent multi-AZ shared storage for your containers
- Note:
 - Amazon S3 cannot be mounted as a file system
-



ECS Service Auto Scaling

- Automatically increase/decrease the desired number of ECS tasks
- Amazon ECS Auto Scaling uses AWS Application Auto Scaling
 - ECS Service Average CPU Utilization
 - ECS Service Average Memory Utilization - Scale on RAM
 - ALB Request Count Per Target - metric coming from the ALB
- **Target Tracking** - scale based on target value for a specific CloudWatch metric
- **Step Scaling** - scale based on a specified CloudWatch Alarm
- **Scheduled Scaling** - scale based on a specified data/time (predictable changes)

- ECS Service Auto Scaling (task level) ≠ EC2 Auto Scaling (EC2 instance level)
- Fargate Auto Scaling is much easier to setup (because Serverless)

EC2 Launch Type - Auto Scaling EC2 Instances

- Accommodate EC2 Service Scaling by adding underlying EC2 Instances
- **Auto Scaling Group Scaling**
 - Scale your ASG based on CPU Utilization
 - Add EC2 instances over time
- **ECS Cluster Capacity Provider**
 - Used to automatically provision and scale the infrastructure for your ECS Tasks
 - Capacity Provider paired with an Auto Scaling Group
 - Add EC2 Instances when you're missing capacity (CPU, RAM...)

Hands on

Cluster configuration

Cluster name
DemoCluster
There can be a maximum of 255 characters. The valid characters are letters (uppercase and lowercase), numbers, hyphens, and underscores.

Networking Info
By default tasks and services run in the default subnets for your default VPC. To use the non-default VPC, specify the VPC and subnets.

VPC
Use a VPC with public and private subnets. By default, VPCs are created for your AWS account. To create a new VPC, go to the [VPC Console](#).
vpc-0a40b157dcd7c81f1
default

Subnets - optional
Select the subnets where your tasks run. We recommend that you use three subnets for production.
Choose subnets
subnet-01a4793dc9adf20bc X eu-west-1c
subnet-09b25e9baf4dae1d7 X eu-west-1a
subnet-069b03c0b1a519b91 X eu-west-1b

Amazon EC2 instances
Manual configurations. Use for large workloads with consistent resource demands.

Auto Scaling group (ASG) [Info](#)
Use Auto Scaling groups to scale the Amazon EC2 instances in the cluster.

Create new ASG

Operating system/Architecture Choose the Windows operating system or Linux architecture for your instance. Amazon Linux 2	EC2 instance type Choose based on the workloads you plan to run on this cluster. t2.micro
Desired capacity Specify the number of instances to launch in your cluster. Minimum 0 Maximum 5	SSH Key pair Create a key pair in the EC2 console, consisting of a private key and a public key, that you use to prove your identity when connecting to an instance. None - unable to SSH
<input type="checkbox"/> External instances using ECS Anywhere Manual configurations. Use to add data center compute.	

► Monitoring - optional [Info](#)
Container Insights is off by default. When you use Container Insights, there is a cost associated with it.

► Tags - optional [Info](#)
Tags help you to identify and organize your clusters.

Cancel **Create**

Cluster DemoCluster has been created successfully.

Amazon Elastic Container Service > Clusters

All Clusters [Info](#)

Clusters (1)		C	Create cluster
DemoCluster No default found Services Tasks No services deployed in cluster. No tasks running in this cluster.		< 1 >	@

EC2 > Auto Scaling groups

Auto Scaling groups (1)							C	Edit	Delete	Create an Auto Scaling group
	Name	Launch template/configuration	Instances	Status	Desired capacity					
<input type="checkbox"/>	Infra-ECS-Cluster	ECSLaunchTemplate_egbArPGLqbbz Ver 0	0	-	0					

Group details	
Edit	
Desired capacity	Auto Scaling group name
0	Infra-ECS-Cluster-0cc5f695-1b06-49a1-9d5f-f1c495ba4dc8-ECSAutoScalingGroup-1PMXWABH1NKSG
Minimum capacity	Date created
0	Thu Apr 07 2022 23:24:42 GMT+0100 (Western European Summer Time)
Maximum capacity	Amazon Resource Name (ARN)
5	arn:aws:autoscaling:eu-west-1:783768293452:autoScalingGroup:80f4be37-38a3-48d3-9919-cdfd5f7f8f39:autoScalingGroupName/Infra-ECS-Cluster-0cc5f695-1b06-49a1-9d5f-f1c495ba4dc8-ECSAutoScalingGroup-1PMXWABH1NKSG

Group size X

Specify the size of the Auto Scaling group by changing the desired capacity. You can also specify minimum and maximum capacity limits. Your desired capacity must be within the limit range.

Desired capacity

Minimum capacity

Maximum capacity

[Cancel](#) [Update](#)

Details		Activity		Automatic scaling		Instance management		Monitoring		Instance refresh									
Instances (1)																			
<input type="button" value="Filter instances"/> Actions ▾ 											< 1 >								
<input type="checkbox"/>	Instance ID	▲	Lifecycle	▼	Instance ty... ▼	Weighted capacity ▼	Launch template/config												
<input type="checkbox"/>	i-04cf214f227bc1e1e		InService		t2.micro	-	ECSLaunchTemplate_egb												

Container instances (1) Info						
<input type="checkbox"/>	Container instance	Status	Type	Instance ID	Availability zo...	Running tasks... CPU available
<input type="checkbox"/>	f973104dc742485...	Active	EC2	i-04cf214f227b...	eu-west-1b	0 1024

Configure task definition and containers

New ECS Experience [Tell us what you think](#)

Amazon Elastic Container Service

- Clusters
- Task definitions**
- Account settings

Amazon ECR

- Repositories

Amazon Elastic Container Service > Task definitions

Task definitions (2) Info		
<input type="checkbox"/>	Task definition	Status of last revision
<input type="radio"/>	demo-nginx	INACTIVE
<input type="radio"/>	first-run-task-definition	INACTIVE

[Create new task definition](#)

Configure task definition and containers

Task definition configuration

Task definition family [Info](#)

Specify a unique task definition family name.

Up to 255 letters (uppercase and lowercase), numbers, hyphens, and underscores are allowed.

Container - 1 [Info](#)

[Essential container](#) [Remove](#)

Container details

Specify a name, container image, and whether the container should be marked as essential. Each task definition must have at least one essential container.

Name	Image URI	Essential container
nginxdemos-hello	nginxdemos/hello	Yes

Port mappings [Info](#)

Add port mappings to allow the container to access ports on the host to send or receive traffic.

Container port	Protocol	Remove
80	TCP	Remove

[Add more port mappings](#)

▼ Environment variables - optional [Info](#)

Add individually

Add a key-value pair to specify an environment variable.

▼ Environment

Specify the infrastructure requirements for the task definition.

App environment | [Info](#)

Specify the infrastructure for the task definition.

[Add an option](#) ▾

AWS Fargate (serverless) [X](#)

Operating system/Architecture | [Info](#)

Linux/X86_64 ▾

Task size | [Info](#)

Specify the amount of CPU and memory to reserve for your task.

CPU

Memory

1 vCPU ▾

3 GB ▾

► Container size - optional [Info](#)

▼ Task roles. network mode - conditional

[awslogs-region](#) [Value](#) [eu-west-1](#)

[awslogs-stream-prefix](#) [Value](#) [ecs](#)

[awslogs-create-group](#) [Value](#) [true](#) [Remove](#)

[Add](#)

Use trace collection [Info](#)

Amazon ECS creates an AWS Distro for OpenTelemetry sidecar to route traces from your application to AWS X-Ray. See pricing information on [AWS X-Ray](#).

Use metric collection [Info](#) [Preview](#)

Amazon ECS creates an AWS Distro for OpenTelemetry sidecar to route custom container and application metrics to Amazon CloudWatch or Amazon Managed Service for Prometheus.

► Tags - optional [Info](#)

Tags help you to identify and organize your task definitions.

[Cancel](#)

[Previous](#)

[Next](#)

Task definition configuration

Task definition family
nginxdemos-hello

Container - 1

Container details

Name nginxdemos-hello	Image URI nginxdemos/hello	Essential container Yes
--------------------------	-------------------------------	----------------------------

Port mappings

Host port:container port/protocol
-:80/tcp

► Environment variables - optional

► Container size - optional

Amazon CloudWatch

awslogs-group	Value	/ecs/nginxdemos-hello
awslogs-region	Value	eu-west-1
awslogs-stream-prefix	Value	ecs
awslogs-create-group	Value	true

Trace collection

Off

Metric collection

Off

Tags (1)

Key	Value
ecs:taskDefinition:createdFrom	ecs-console-v2

Create

EC2 > Security Groups > Create security group

Create security group Info

A security group acts as a virtual firewall for your instance to control inbound and outbound traffic. To create a new security group, complete the fields below.

Basic details

Security group name Info
alb-ecs-sg

Name cannot be edited after creation.

Description Info
ALB for ECS SG

VPC Info
vpc-0a40b157dc7c81f1

Inbound rules Info

Type <small>Info</small>	Protocol <small>Info</small>	Port range <small>Info</small>	Source <small>Info</small>	Description - optional <small>Info</small>
Custom TCP	TCP	80	Anywhere <input type="button" value="Search"/> 0.0.0.0/0 <input type="button" value="X"/>	<input type="text"/> Delete
Custom TCP	TCP	80	Anywhere <input type="button" value="Search"/> ::/0 <input type="button" value="X"/>	<input type="text"/> Delete

[Add rule](#)

Outbound rules Info

Type <small>Info</small>	Protocol <small>Info</small>	Port range <small>Info</small>	Destination <small>Info</small>	Description - optional <small>Info</small>
All traffic	All	All	Custom <input type="button" value="Search"/> 0.0.0.0/0 <input type="button" value="X"/>	<input type="text"/> Delete

[Add rule](#)

Tags - optional

A tag is a label that you assign to an AWS resource. Each tag consists of a key and an optional value. You can use tags to search and filter your resources or track your AWS costs.

No tags associated with the resource.

[Add new tag](#)
You can add up to 50 more tag

[Cancel](#) [Create security group](#)

Create security group Info

A security group acts as a virtual firewall for your instance to control inbound and outbound traffic. To create a new security group, complete the fields below.

Basic details

Security group name <small>Info</small>	<input type="text" value="nginx-demo-sg"/>
Name cannot be edited after creation.	
Description <small>Info</small>	<input type="text" value="Allows SSH access to developers"/>
VPC <small>Info</small>	<input type="text" value="vpc-0a40b157dcd7c81f1"/> <input type="button" value="X"/>

Inbound rules Info

Type <small>Info</small>	Protocol <small>Info</small>	Port range <small>Info</small>	Source <small>Info</small>	Description - optional <small>Info</small>
All TCP	TCP	0 - 65535	Custom <input type="button" value="Search"/> sg-02ad0dad0206f9b37 <input type="button" value="X"/>	Allow traffic from the ALB <input type="text"/> Delete

[Add rule](#)

Outbound rules [Info](#)

Type	Info	Protocol	Info	Port range	Info	Destination	Info	Description - optional	Info
All traffic	▼	All	All	Custom	▼	<input type="text" value="Custom"/>	<input type="text" value=""/>	0.0.0.0/0	X
Delete									

[Add rule](#)

Tags - optional
A tag is a label that you assign to an AWS resource. Each tag consists of a key and an optional value. You can use tags to search and filter your resources or track your AWS costs.

No tags associated with the resource.

[Add new tag](#)
You can add up to 50 more tag

[Cancel](#) [Create security group](#)

Task definition successfully created
nginxdemos-hello has been successfully created. You can use this task definition to deploy a service or run a task.

[Amazon Elastic Container Service](#) > Clusters

All Clusters [Info](#)

Clusters (1)		Create cluster				
<input type="text" value="Search clusters"/>		C 1 G				
DemoCluster <u>No default found</u> <table border="1"> <thead> <tr> <th>Services</th> <th>Tasks</th> </tr> </thead> <tbody> <tr> <td>No services deployed in cluster.</td> <td>No tasks running in this cluster.</td> </tr> </tbody> </table>			Services	Tasks	No services deployed in cluster.	No tasks running in this cluster.
Services	Tasks					
No services deployed in cluster.	No tasks running in this cluster.					

Services	Tasks	Infrastructure	Metrics	Tags																		
<h3>Services (0)</h3> <table border="1"> <thead> <tr> <th>C</th> <th>Edit</th> <th>Delete</th> <th>Deploy</th> </tr> </thead> <tbody> <tr> <td colspan="4"><input type="text" value="Filter services by value"/></td> </tr> <tr> <td>Service name ▾</td> <td>ARN</td> <td>Status</td> <td>Deployments and tasks</td> <td>Task definiti... ▾</td> </tr> <tr> <td>Revision</td> <td>Laun...</td> <td colspan="3"></td> </tr> </tbody> </table> <p>No services No services to display. Deploy</p>					C	Edit	Delete	Deploy	<input type="text" value="Filter services by value"/>				Service name ▾	ARN	Status	Deployments and tasks	Task definiti... ▾	Revision	Laun...			
C	Edit	Delete	Deploy																			
<input type="text" value="Filter services by value"/>																						
Service name ▾	ARN	Status	Deployments and tasks	Task definiti... ▾																		
Revision	Laun...																					

Application type [Info](#)

Specify what type of application you want to run.

 Service

Launch a group of tasks handling a long-running computing work that can be stopped and restarted. For example, a web application.

 Task

Launch a standalone task that runs and terminates. For example, a batch job.

Task definitionSelect an existing task definition. To create a new task definition, go to [Task definitions](#). **Specify revision manually**

Manually input revision instead of choosing from the 100 most recent revisions for the selected task definition family.

Family**Revision**

nginxdemos-hello

▼

1

▼

Service name

Assign a unique name for this service.

Desired tasks

Specify the number of tasks to launch.

1

▼ Load balancing - optional**Load balancer type** [Info](#)

Configure a load balancer to distribute incoming traffic across the tasks running in your service.

Application Load Balancer

▼

Application Load Balancer

Specify whether to create a new load balancer or choose an existing one.

- Create a new load balancer**
- Use an existing load balancer**

Load balancer name

Assign a unique name for the load balancer.

DemoALBForECS

Listener [Info](#)

Specify the port and protocol that the load balancer will listen for connection requests on.

Port**Protocol**

80

HTTP

▼

Target group [Info](#)

Create a target group that the load balancer will use to route requests to the tasks in your service.

Target group name**Protocol**

nginx-ecs

HTTP

▼

Health check path [Info](#)

/

Health check grace period [Info](#)

20

seconds

Subnets
Choose the subnets within the VPC that the task scheduler should consider for placement.

Choose subnets ▾

- subnet-01a4793dc9adf20bc X
eu-west-1c
- subnet-09b25e9baf4dae1d7 X
eu-west-1a
- subnet-069b03c0b1a519b91 X
eu-west-1b

Security group [Info](#)
Choose an existing security group or create a new security group.

Use an existing security group

Create a new security group

Security group name
Choose an existing security group.

sg-0284edb2b9ef37ffe X
nginx-demo-sg

Public IP [Info](#)
Choose whether to auto-assign a public IP to the task's elastic network interface (ENI).

Enabled

[Create Load Balancer](#) [Actions ▾](#)

Filter by tags and attributes or search by keyword 1 to 1 of 1

Name	DNS name	State	VPC ID	Availability Zones
DemoALBForECS	DemoALBForECS-1739327...	Active	vpc-0a40b157dcd7c81f1	eu-west-1c, eu-west-1b...

IPv4 address: Assigned by AWS

[Edit subnets](#)

Hosted zone Z32O12XQLNTSW2
Creation time April 8, 2022 at 12:17:09 AM UTC+1

Security

Security groups sg-02ad0dadcc206f9b37, alb-ecs-sg
ALB for ECS SG

[Edit security groups](#)

Attributes

Load balancer: DemoALBForECS

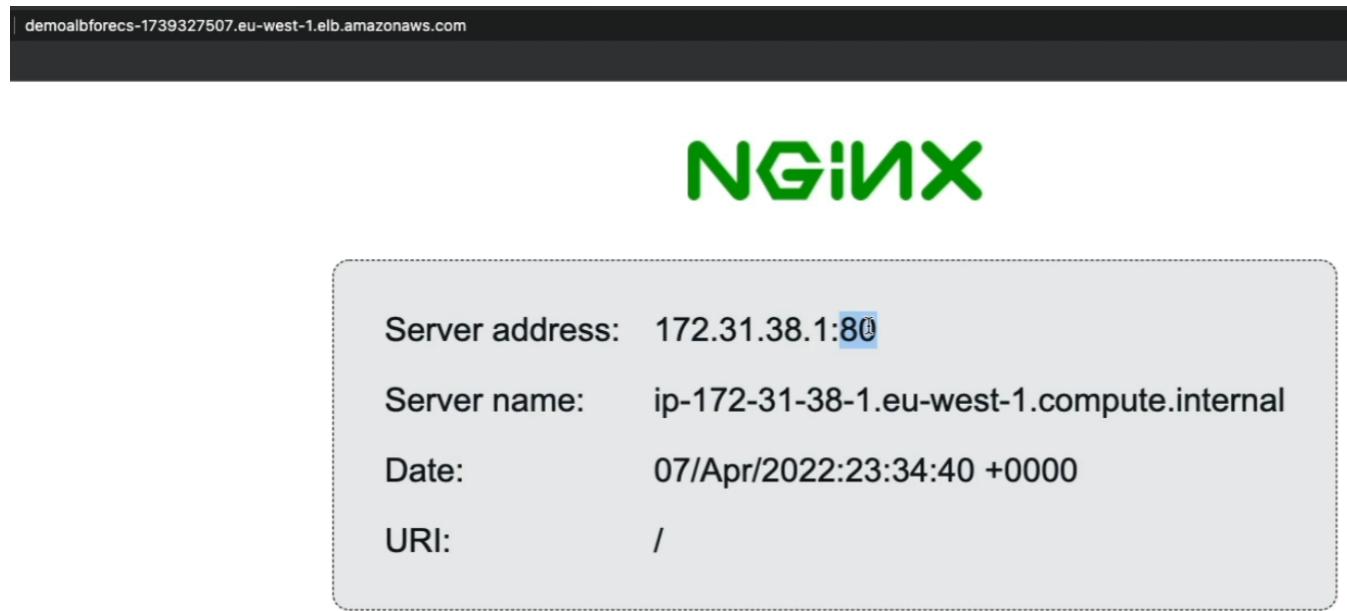
Description [Listeners](#) Monitoring Integrated services Tags

Listeners listen for connection requests using their protocol and port. You can add, remove, or update listeners and listener rules.

To view and edit listener attributes, select the listener and choose Edit.

Add listener [Edit](#) [Delete](#)

Listener ID	Security policy	SSL Certificate	Rules
HTTP : 80	N/A	N/A	Default: forwarding to nginx-ecs View/edit rules
arn...e570b9a7cf819ca7			



Deployment configuration

Task definition
 Select an existing task definition. To create a new task definition, go to [Task definitions](#).

Specify revision manually
 Manually input revision instead of choosing from the 100 most recent revisions for the selected task definition family.

Family	Revision
nginxdemos-hello	1

Desired tasks
 Specify the number of tasks to launch.

► Deployment options

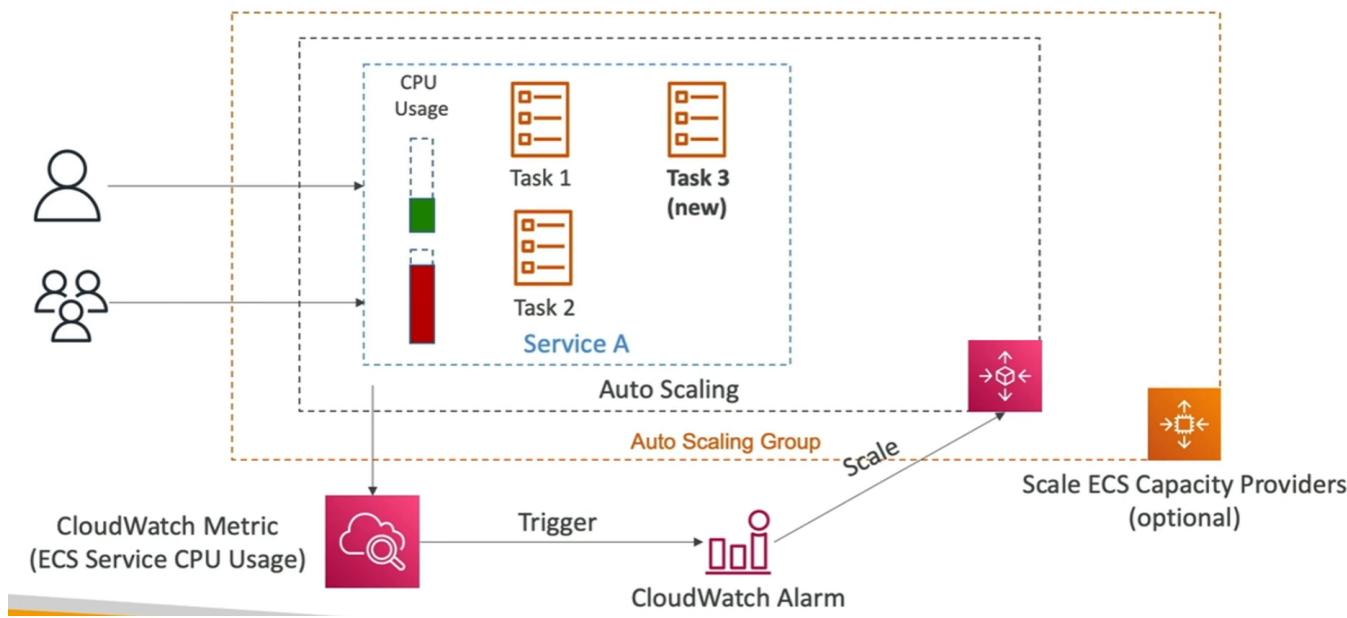
► Tags - optional
 Tags help you to identify and organize your clusters.

Cancel **Update**

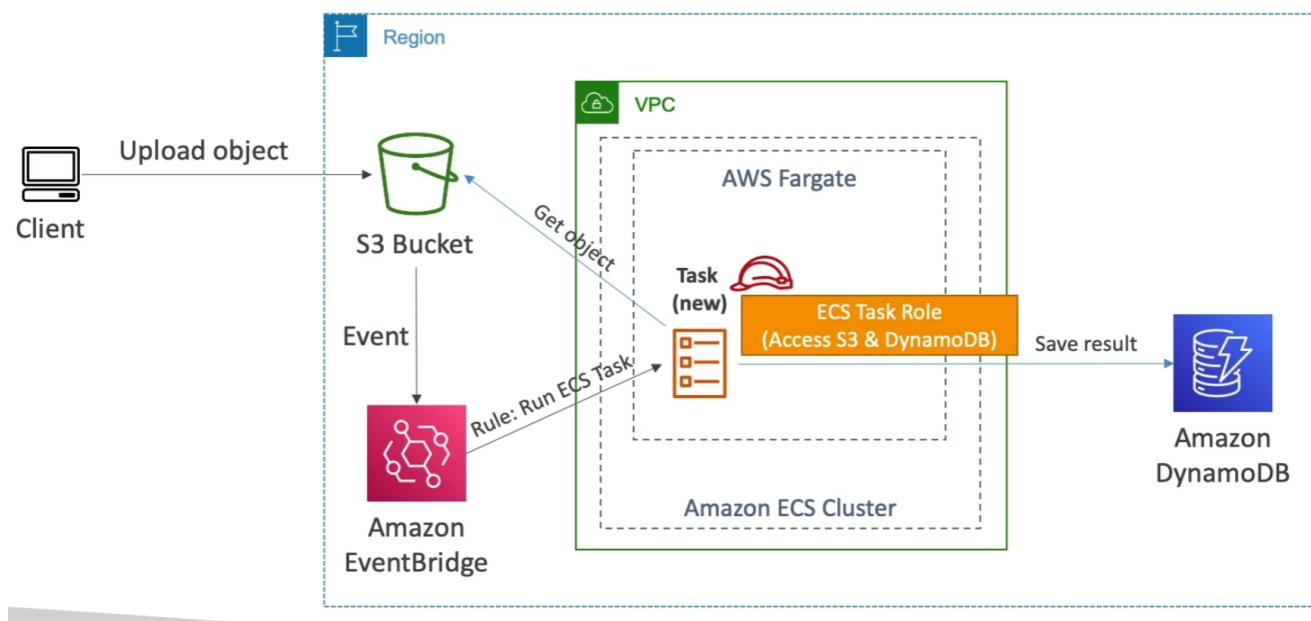
Tasks (4)

Task	Last status	Task definiti...	Revisi...	Health stat...	Started...	Container instan...	Launch type
6db6b...	Running	nginxdemos-hello	1	Unknown	15 minute...	-	FARGATE
Geb49...	Provisioning	nginxdemos-hello	1	Unknown	-	-	FARGATE
72963...	Pending	nginxdemos-hello	1	Unknown	-	-	FARGATE
759f92...	Pending	nginxdemos-hello	1	Unknown	-	-	FARGATE

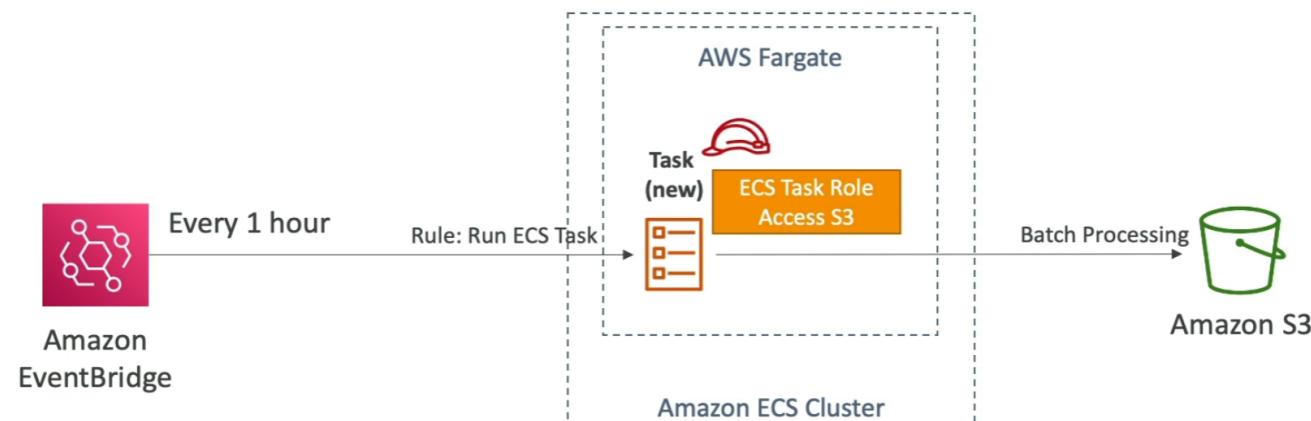
ECS Scaling - Service CPU Usage Example



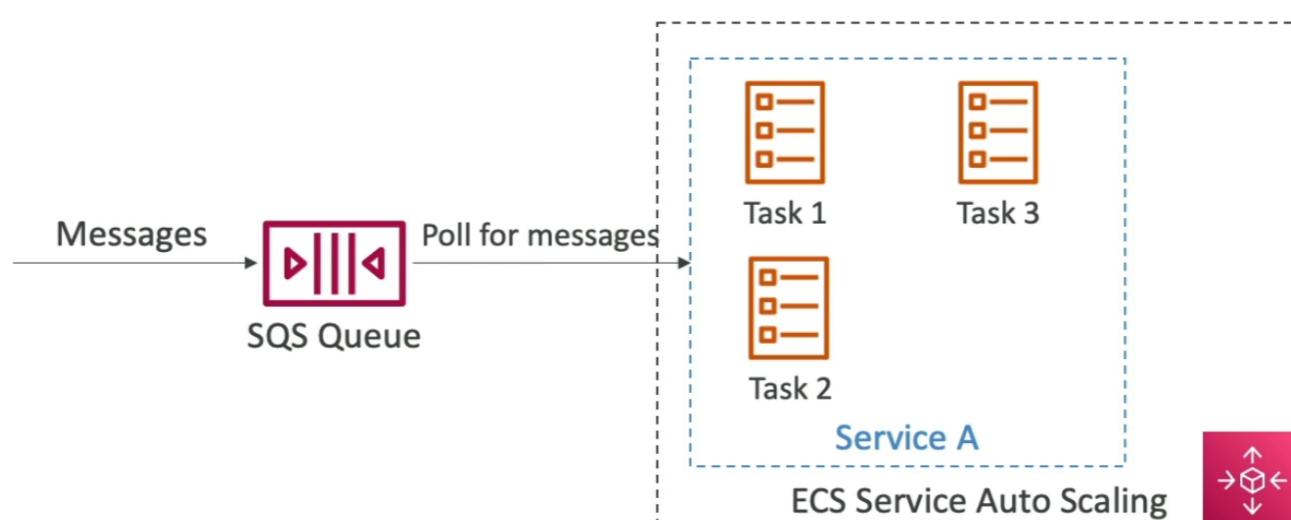
ECS tasks invoked by Event Bridge



ECS tasks invoked by Event Bridge Schedule

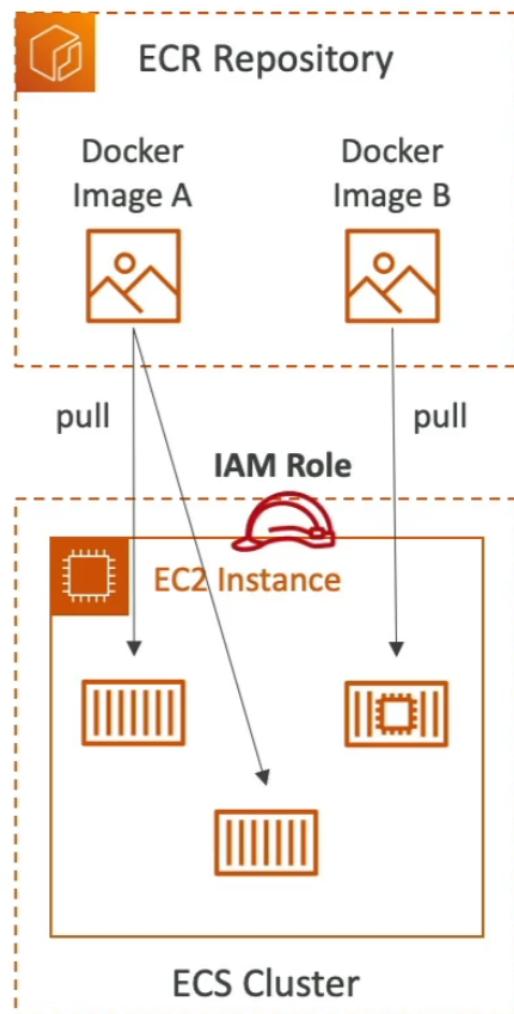


ECS - SQS Queue Example



Amazon ECR

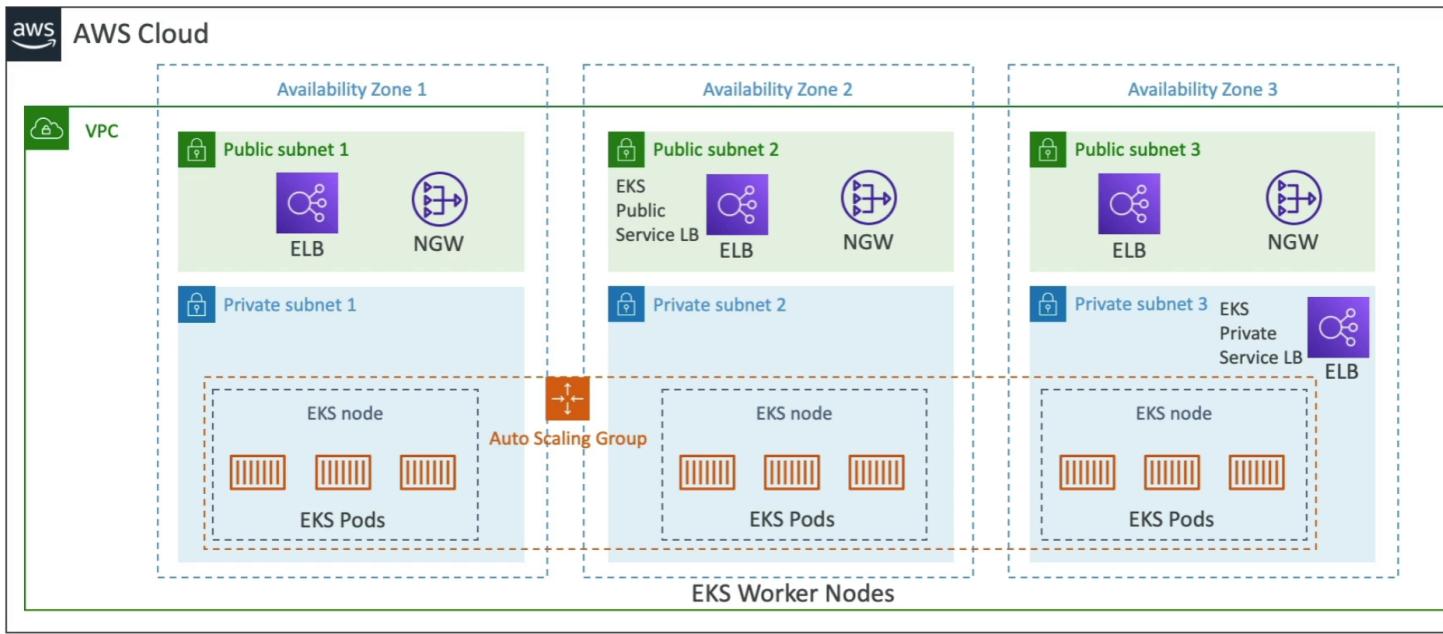
- ECR = Elastic Container Registry
- Store and manage Docker images on AWS
- Private and Public repository (Amazon ECR Public Gallery)
- Fully integrated with ECS, backed by Amazon S3
- Access is controlled through IAM (permission errors => policy)
- Supports image vulnerability scanning, versioning, image tags, image lifecycle, ...



Amazon EKS Overview

- Amazon EKS = Amazon Elastic **Kubernetes** Service
- It is a way to launch **managed Kubernetes clusters on AWS**
- Kubernetes is an open-source system for automatic deployment, scaling and management of containerized (usually Docker) application
- It's an alternative to ECS, similar goal but different API
- EKS supports **EC2** if you want to deploy worker nodes or **Fargate** to deploy serverless containers
- **Use case:** if your company is already using Kubernetes on-premises or in another cloud, and wants to migrate to AWS using Kubernetes
- Kubernetes is cloud-agnostic (can be used in any cloud - Azure, GCP..)

Amazon EKS - Diagram



Amazon EKS - Node Types

- **Managed Node Groups**
 - Create and manages Nodes (ECs instances) for you
 - Nodes are part of an ASG managed by EKS
 - Supports On-Demand or Spot Instances
- **Self-Managed Nodes**
 - Nodes created by you and registered to the EKS cluster and managed by an ASG
 - You can use prebuilt AMI - Amazon EKS Optimized AMI
 - Supports On-Demand or Spot Instances
- **AWS Fargate**
 - No maintenance required; no nodes managed

Amazon EKS - Data Volumes

- Need to specify **StorageClass** manifest on your EKS cluster
- Leverages a **Container Storage Interface (CSI)** compliant driver
- **Support for ...**
 - Amazon EBS
 - Amazon EFS (work with Fargate)
 - Amazon FSx for Lustre
 - Amazon FSx for NetApp ONTAP

Hands on

Step 1
Configure cluster

Step 2
Specify networking

Step 3
Configure logging

Step 4
Review and create

Configure cluster

Cluster configuration Info

Name
Enter a unique name for this cluster. This property cannot be changed after the cluster is created.

The cluster name should begin with letter or digit and can have any of the following characters: the set of Unicode letters, digits, hyphens and underscores. Maximum length of 100.

Kubernetes version Info
Select the Kubernetes version for this cluster.

Cluster service role Info
Select the IAM role to allow the Kubernetes control plane to manage AWS resources on your behalf. This property cannot be changed after the cluster is created. To create a new role, follow the instructions in the [Amazon EKS User Guide](#).

Identity and Access Management (IAM)

New! Securely access AWS services from your data center with IAM Roles Anywhere. [Learn more](#)

IAM > Roles

Roles (45) Info

An IAM role is an identity you can create that has specific permissions with credentials that are valid for short durations. Roles can be assumed by entities that you trust.

Role name	Trusted entities
AmazonEC2RoleForCloudWatch	AWS Service: ec2
AmazonSSMRoleForInstancesQuickSetup	AWS Service: ec2
amplify-login-lambda-e258514f	AWS Service: lambda
aws-ec2-spot-fleet-tagging-role	AWS Service: spotfleet
aws-elasticbeanstalk-ec2-role	AWS Service: ec2
aws-elasticbeanstalk-service-role	AWS Service: elasticbeanstalk
AWS-QuickSetup-HostMgmtRole-eu-west-1-qeizg	AWS Service: ssm

IAM > Roles > Create role

Step 1
Select trusted entity

Step 2
Add permissions

Step 3
Name, review, and create

Select trusted entity

Trusted entity type

- AWS service**
Allow AWS services like EC2, Lambda, or others to perform actions in this account.
- AWS account**
Allow entities in other AWS accounts belonging to you or a 3rd party to perform actions in this account.
- Web identity**
Allows users federated by the specified external web identity provider to assume this role to perform actions in this account.
- SAML 2.0 federation**
Allow users federated with SAML 2.0 from a corporate directory to perform actions in this account.
- Custom trust policy**
Create a custom trust policy to enable others to perform actions in this account.

Use cases for other AWS services:

EKS

- EKS**
Allows EKS to manage clusters on your behalf.
- EKS - Cluster**
Allows access to other AWS service resources that are required to operate clusters managed by EKS.
- EKS - Nodegroup**
Allow EKS to manage nodegroups on your behalf.
- EKS - Fargate pod**
Allows access to other AWS service resources that are required to run Amazon EKS pods on AWS Fargate.
- EKS - Fargate profile**
Allows EKS to run Fargate tasks.
- EKS - Connector**
Allows access to other AWS service resources that are required to connect to external clusters
- EKS Local - Outpost**
Allows Amazon EKS Local to call AWS services on your behalf.

[Cancel](#) [Next](#)

Create IAM Role

Step 2
Add permissions

Role details

Step 3
Name, review, and create

Role name
Enter a meaningful name to identify this role.

Maximum 64 characters. Use alphanumeric and '+,-,@-' characters.

Description
Add a short explanation for this role.

Maximum 1000 characters. Use alphanumeric and '+,-,@-' characters.

Step 1: Select trusted entities

Edit

```

1  [
2   "Version": "2012-10-17",
3   "Statement": [
4     {
5       "Effect": "Allow",
6       "Principal": {
7         "Service": [
8           "eks.amazonaws.com"
9         ]
10      }
11    }
12  ]
13 }
```

Step 2: Add permissions

Edit

Permissions policy summary

Policy name	Type	Attached as
AmazonEKSClusterPolicy	AWS managed	Permissions policy

Tags

Add tags (Optional)

Tags are key-value pairs that you can add to AWS resources to help identify, organize, or search for resources.

No tags associated with the resource.

[Add tag](#)

You can add up to 50 more tags.

[Cancel](#) [Previous](#) [Create role](#)

Select created IAM Role

The cluster name should begin with letter or digit and can have any of the following characters: the set of Unicode letters, digits, hyphens and underscores. Maximum length of 100.

Kubernetes version [Info](#)
Select the Kubernetes version for this cluster.

Cluster service role [Info](#)
Select the IAM role to allow the Kubernetes control plane to manage AWS resources on your behalf. This property cannot be changed after the cluster is created. To create a new role, follow the instructions in the [Amazon EKS User Guide](#).

[C](#)

Networking Setting

Networking Info

These properties cannot be changed after the cluster is created.

VPC Info

Select a VPC to use for your EKS cluster resources. To create a new VPC, go to the [VPC console](#).

vpc-0a40b157dc7c81f1 | Default



Subnets Info

Choose the subnets in your VPC where the control plane may place elastic network interfaces (ENIs) to facilitate communication with your cluster. To create a new subnet, go to the corresponding page in the [VPC console](#).

Select subnets



subnet-01a4793dc9adf20bc X

subnet-09b25e9baf4dae1d7 X

subnet-069b03c0b1a519b91 X



Security groups Info

Choose the security groups to apply to the EKS-managed Elastic Network Interfaces that are created in your worker node subnets. To create a new security group, go to the corresponding page in the [VPC console](#).

Select security groups



Choose cluster IP address family Info

Specify the IP address type for pods and services in your cluster.

- IPv4
 IPv6

Networking add-ons

Configure add-ons that provide advanced networking functionalities on the cluster.

Amazon VPC CNI Info

Enable pod networking within your cluster.

Version

Select the version for this add-on.

v1.10.1-eksbuild.1



i This add-on will use the IAM role of the node where it runs. You can change this add-on to use IAM roles for service accounts after cluster creation.

CoreDNS Info

Enable service discovery within your cluster.

Version

Select the version for this add-on.

v1.8.7-eksbuild.1



kube-proxy Info

Enable service networking within your cluster.

Version

Configure logging

Control plane logging [Info](#)
Send audit and diagnostic logs from the Amazon EKS control plane to CloudWatch Logs.

API server
Logs pertaining to API requests to the cluster.

Audit
Logs pertaining to cluster access via the Kubernetes API.

Authenticator
Logs pertaining to authentication requests into the cluster.

Controller manager
Logs pertaining to state of cluster controllers.

Scheduler
Logs pertaining to scheduling decisions.

[Cancel](#) [Previous](#) [Next](#)

Step 3: Logging [Edit](#)

Control plane logging

API server off	Audit off	Authenticator off
Controller manager off	Scheduler off	

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

Configure node group

Configure node group [Info](#)
A node group is a group of EC2 instances that supply compute capacity to your Amazon EKS cluster. You can add multiple node groups to your cluster.

Node group configuration
These properties cannot be changed after the node group is created.

Name
Assign a unique name for this node group.

The node group name should begin with letter or digit and can have any of the following characters: the set of Unicode letters, digits, hyphens and underscores. Maximum length of 63.

Node IAM role [Info](#)
Select the IAM role that will be used by the nodes. To create a new role, go to the [IAM console](#).

[C](#)

Info The selected role must not be used by a self-managed node group as this could lead to a service interruption upon managed node group deletion.
[Learn more](#)

Use case

Allow an AWS service like EC2, Lambda, or others to perform actions in this account.

Common use cases

- EC2**
Allows EC2 instances to call AWS services on your behalf.
- Lambda**
Allows Lambda functions to call AWS services on your behalf.

Use cases for other AWS services:

Choose a service to view use case ▾

Cancel

Next

Add permissions

Permissions policies (Selected 1/777)

Choose one or more policies to attach to your new role.

Filter policies by property or policy name and press enter

7 matches < 1 >

"eks" X Clear filters

Policy name	Type	Description
AmazonEKSClusterPolicy	AWS m...	This policy provi
<input checked="" type="checkbox"/> AmazonEKSWorkerNodePolicy	AWS m...	This policy allow
AmazonEKSServicePolicy	AWS m...	This policy allow
AmazonEKS_CNI_Policy	AWS m...	This policy provi
AmazonEKSFargatePodExecutionRolePolicy	AWS m...	Provides access
AmazonEKSLocalOutpostClusterPolicy	AWS m...	This policy provi

Add permissions

Permissions policies (Selected 2/777)

Choose one or more policies to attach to your new role.

Filter policies by property or policy name and press enter

1 match < 1 > ⚙

"AmazonEC2ContainerRegistryReadOnly" X Clear filters

Policy name	Type	Description
<input checked="" type="checkbox"/> AmazonEC2ContainerRegistryReadOnly	AWS m...	Provides read-only acce

► Set permissions boundary - optional

Set a permissions boundary to control the maximum permissions this role can have. This is not a common setting, but you can use it to delegate permission management to others.

Step 2: Add permissions

Edit

Permissions policy summary		
Policy name	Type	Attached as
AmazonEC2ContainerRegistryReadOnly	AWS managed	Permissions policy
AmazonEKSWorkerNodePolicy	AWS managed	Permissions policy

Tags

Add tags (Optional)

Tags are key-value pairs that you can add to AWS resources to help identify, organize, or search for resources.

No tags associated with the resource.

Add tag

You can add up to 50 more tags.

Cancel Previous Create role

Node group configuration

These properties cannot be changed after the node group is created.

Name

Assign a unique name for this node group.

DemoNodeGroup

The node group name should begin with letter or digit and can have any of the following characters: the set of Unicode letters, digits, hyphens and underscores. Maximum length of 63.

Node IAM role [Info](#)

Select the IAM role that will be used by the nodes. To create a new role, go to the [IAM console](#).

AmazonEKSNodeRole



i The selected role must not be used by a self-managed node group as this could lead to a service interruption upon managed node group deletion.

[Learn more](#)

Node group compute configuration

These properties cannot be changed after the node group is created.

AMI type [Info](#)

Select the EKS-optimized Amazon Machine Image for nodes.

Amazon Linux 2 (AL2_x86_64)



Capacity type

Select the capacity purchase option for this node group.

On-Demand



Instance types [Info](#)

Select instance types you prefer for this node group.

Select



t3.medium



vCPU: Up to 2 vCPUs memory: 4.0 GiB

Disk size

Select the size of the attached EBS volume for each node.

20



GiB

Node group scaling configuration

Desired size
Set the desired number of nodes that the group should launch with initially.

2 nodes

Minimum size
Set the minimum number of nodes that the group can scale in to.

2 nodes

Maximum size
Set the maximum number of nodes that the group can scale out to.

2 nodes

Specify networking

Node group network configuration

These properties cannot be changed after the node group is created.

Subnets [Info](#)
Specify the subnets in your VPC where your nodes will run. To create a new subnet, go to the corresponding page in the [VPC console](#).

Select subnets ▾ [C](#)

subnet-01a4793dc9adf20bc X subnet-09b25e9baf4dae1d7 X
subnet-069b03c0b1a519b91 X

Configure SSH access to nodes [Info](#)

Cancel Previous Next

Step 3: Networking

Node group network configuration

Subnets

subnet-01a4793dc9adf20bc	Configure SSH access to nodes
subnet-09b25e9baf4dae1d7	off
subnet-069b03c0b1a519b91	

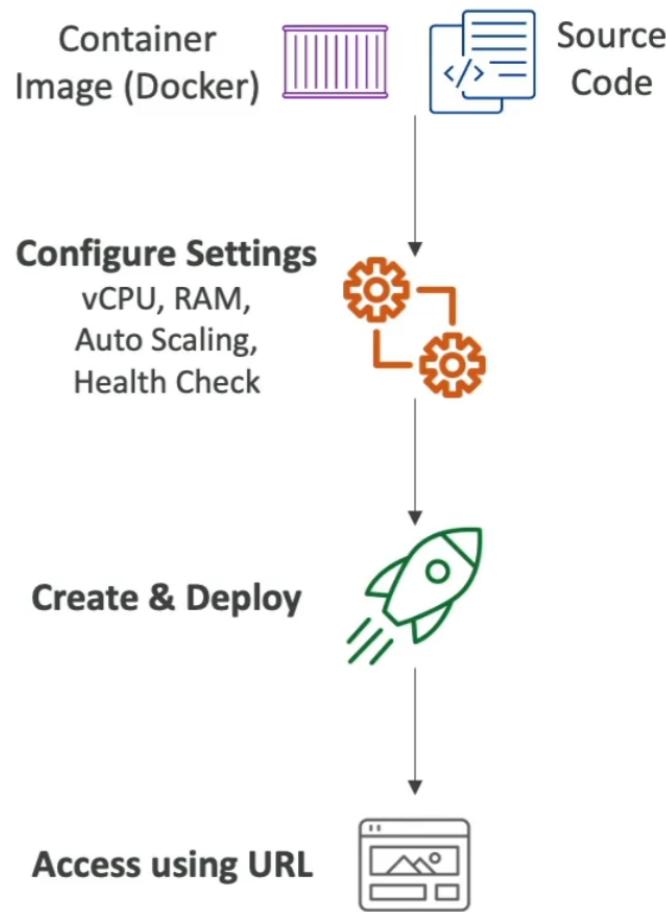
Cancel Previous Create

Node groups (1) Info		Edit	Delete	Add node group
Group name	Desired size	AMI release version	Launch template	Status
DemoNodeGroup	2	1.22.12-20220824	-	Creating

AWS App Runner

- Fully managed service that makes it easy to deploy web applications and APIs at scale
- No infrastructure experience required
- Start with your source code or container image
- Automatically builds and deploy the web app
- Automatic scaling, high available, load balancer, encryption
- VPC access support
- Connect to database, cache, and message queue services

- Use cases: web apps, APIs, microservices, rapid production deployments



Hands on

Repository type

Container registry
 Deploy your service from a container image stored in a container registry.

Source code repository
 Deploy your service from code hosted in a source code repository.

Provider

Amazon ECR

Amazon ECR Public

Container image URI
 Enter a URI to an image you can access, or browse images in your Amazon ECR account.

Deployment settings

Deployment trigger

Manual
 Start each deployment yourself using the App Runner console or AWS CLI.

Automatic
 App Runner monitors your registry and deploys a new version of your service for each image push.

Cancel
Next

Service settings

Service name
DemoHTTP

Virtual CPU & memory
1 vCPU ▾ 2 GB ▾

Environment variables — *optional*
Key-value pairs that you can use to store custom configuration values.
No environment variables have been configured.

Add environment variable

Port
Your service uses this TCP port.
8080 ▾

► Additional configuration

► Networking Info
Configure the way your service communicates with other applications, services, and resources.

► Observability
Configure observability tooling.

▼ Tags Info
Use tags to search and filter your resources, track your AWS costs, and control access permissions.

Tags — *optional*
A tag is a key-value pair that you assign to an AWS resource.
No tags associated with the resource.

Add new tag

You can add 50 more tags.

Cancel Previous Next

⌚ Create service succeeded.

App Runner > Services > DemoHTTP

DemoHTTP Info

Actions ▾ C Deploy

Service overview

Status Running	Service ARN arn:aws:apprunner:eu-west-1:783768293452:service/DemoHTTP /4ba93e36c0d1425b90c2f28c2e6cc3f4
Default domain https://2mdtusr9rz.eu-west-1.awsapprunner.com	Source public.ecr.aws/docker/library/httpd:latest

Logs Activity Metrics Observability Configuration Custom domains

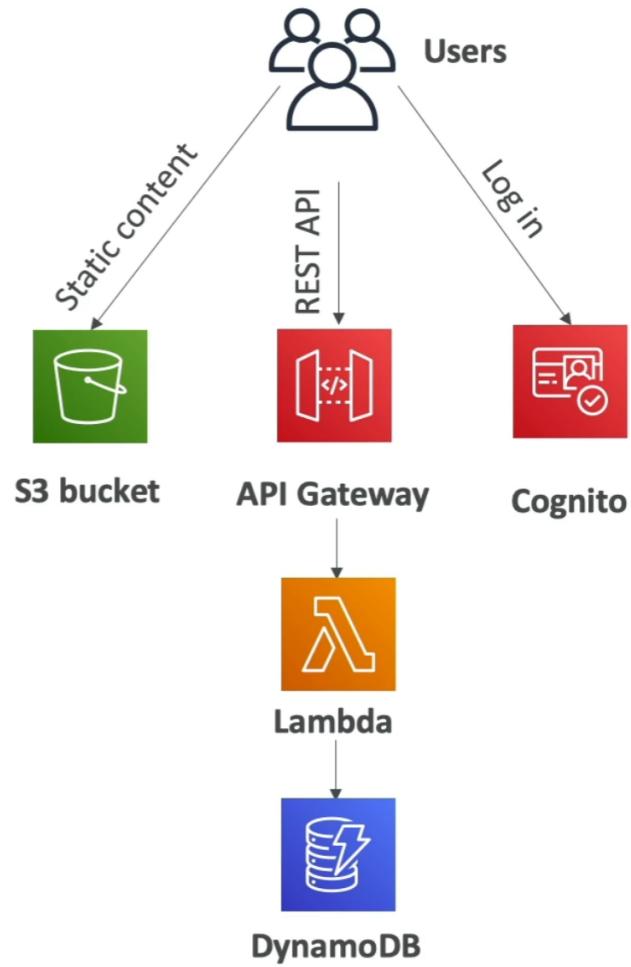
Event log Info

Search

What's serverless?

- Serverless is a new paradigm in which the developers don't have to manage servers anymore
- They just deploy code
- They just deploy ... functions !
- Initially ... Serverless == FaaS (Function as a Service)
- Serverless was pioneered by AWS Lambda but now also includes anything that's managed: "databases,messaging, storage,etc."
- Serverless does not mean there are no servers ...
 - it means you just don't manage / provision /see them

Serverless in AWS



- AWS Lambda
- DynamoDB
- AWS Cognito
- AWS API Gateway
- Amazon S3
- AWS SNS & SQS
- AWS Kinesis Data Firehose
- Aurora Serverless
- Step Functions
- Fargate

Why AWS Lambda

- Amazon EC2
 - Virtual Servers in the Cloud
 - Limited by RAM and CPU
 - Continuously Running
 - Scaling means intervention to add / remove servers
- Amazon Lambda
 - Virtual functions - no servers to manage!
 - Limited by time - short executions
 - Run on-demand
 - Scaling is automated

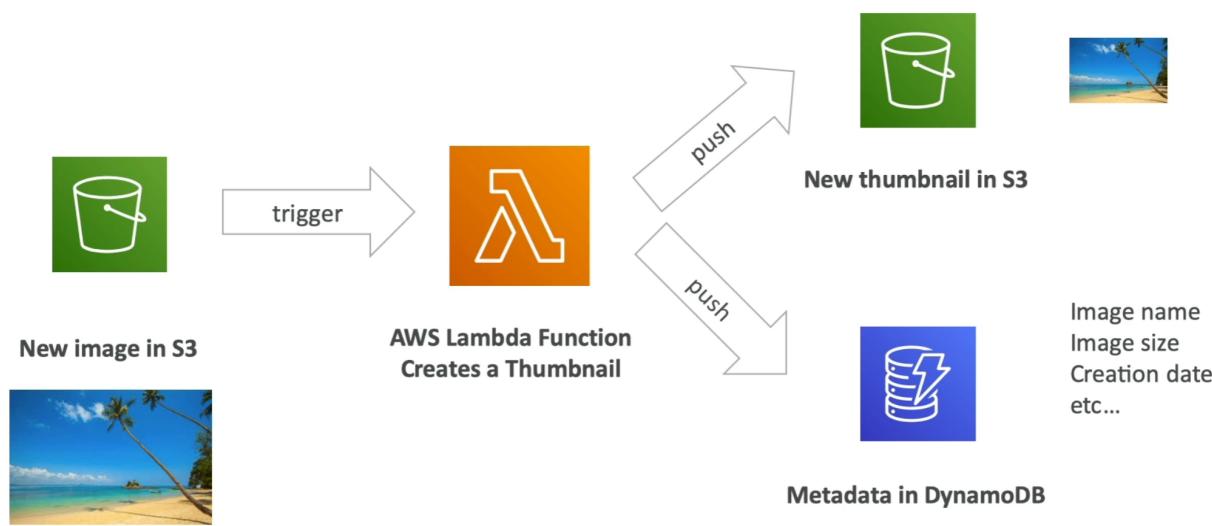
Benefits of AWS Lambda

- Easy Pricing
 - Pay per request and compute time
 - Free tier of 1,000,000 AWS Lambda requests and 400,000 GBs of compute time

- Integrated with the whole AWS suite of services
- Integrated with many programming languages
- Easy monitoring through AWS CloudWatch
- Easy to get more resources per functions (up to 10GB of RAM!)
- Increasing RAM will also improve CPU and network

AWS Lambda Integrations Main ones

- Examole: Serverless Thumbnail creation



- Example: Serverless CRON Job



Hands on

AWS Lambda

Resources for Europe (Ireland)

Lambda function(s)	Code storage	Full account concurrency	Unreserved account concurrency
0	0 byte (0% of 75.0 GB)	1000	1000

Account-level metrics

The charts below show metrics across all your Lambda functions in this AWS Region.

Error count and success rate	Throttles	Invocations
1 to 100 No data available. 0.5 adjusting the dashboard time range.	1 No data available. 0.5 adjusting the dashboard time range.	1 No data available. 0.5 adjusting the dashboard time range.

How it works

Mobile / IoT backends: Represented by a smartphone icon.

Streaming analytics: Represented by a valve icon.

Data processing: Represented by a camera icon.

AWS Lambda: Represented by an orange Lambda icon.

A callout box shows log output from Lambda:

```

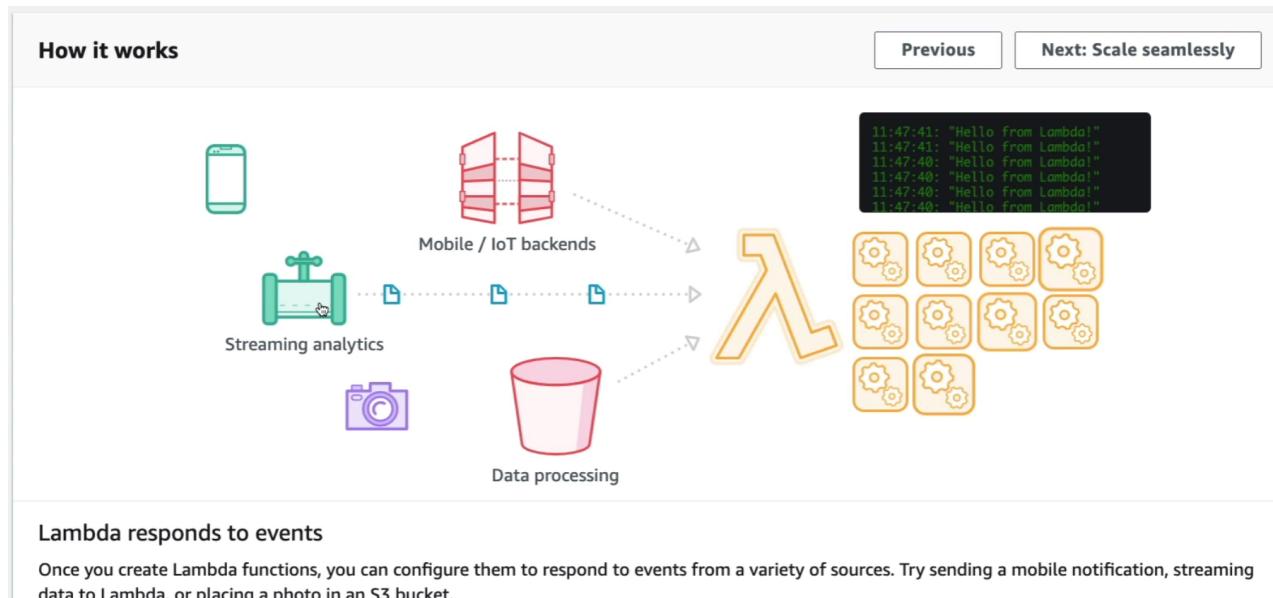
11:47:25: "Hello from Lambda!"

```

Previous | **Next: Scale seamlessly**

Lambda responds to events

Once you create Lambda functions, you can configure them to respond to events from a variety of sources. Try sending a mobile notification, streaming data to Lambda, or placing a photo in an S3 bucket.



Lambda responds to events

Once you create Lambda functions, you can configure them to respond to events from a variety of sources. Try sending a mobile notification, streaming data to Lambda, or placing a photo in an S3 bucket.

Lambda > Functions > Create function

Create function Info

Choose one of the following options to create your function.

- Author from scratch Start with a simple Hello World example.
- Use a blueprint Build a Lambda application from sample code and configuration presets for common use cases.
- Container image Select a container image to deploy for your function.
- Browse serverless app repository Deploy a sample Lambda application from the AWS Serverless Application Repository.

Blueprints Info

Filter by tags and attributes or Search by keyword 1 match

Blueprint name: hello-world-python X Clear filters

Export

Lambda > Functions > Create function > Configure blueprint hello-world-python

Basic information Info

Function name

Execution role Choose a role that defines the permissions of your function. To create a custom role, go to the [IAM console](#).

- Create a new role with basic Lambda permissions
- Use an existing role
- Create a new role from AWS policy templates

ⓘ Role creation might take a few minutes. Please do not delete the role or edit the trust or permissions policies in this role.

Lambda will create an execution role named demo-lambda-role-eaek950w, with permission to upload logs to Amazon CloudWatch Logs.

Lambda function code

Code is preconfigured by the chosen blueprint. You can configure it after you create the function. [Learn more](#) about deploying Lambda functions.

Runtime Python 3.7

```

1 import json
2
3 print('Loading function')
4
5
6 def lambda_handler(event, context):
7     #print("Received event: " + json.dumps(event, indent=2))
8     print("value1 = " + event['key1'])
9     print("value2 = " + event['key2'])
10    print("value3 = " + event['key3'])
11    return event['key1'] # Echo back the first key value
12    #raise Exception('Something went wrong')
13

```

```

1 import json
2
3 print('Loading function')
4
5
6 def lambda_handler(event, context):
7     #print("Received event: " + json.dumps(event, indent=2))
8     print("value1 = " + event["key1"])
9     print("value2 = " + event["key2"])
10    print("value3 = " + event["key3"])
11    return event['key1'] # Echo back the first key value
12    #raise Exception('Something went wrong')
13

```

Configure test event

A function can have up to 10 test events. The events are persisted so you can switch to another computer or web browser and test your function with the same events.

Create new test event
 Edit saved test events

Event template
 hello-world ▾

Event name
 DemoEvent

```

1 {
2     "key1": "value1",
3     "key2": "value2",
4     "key3": "value3"
5 }

```

AWS Lambda Limits to Know - per region

- Execution:
 - Memory allocation: 128MB - 10GB (1 MB increments)
 - Maximum execution time: 900 seconds (15 minutes)
 - Environment variables (4KB)
 - Disk capacity in the "function container" (in /tmp): 512MB to 10GB
 - Concurrency executions: 1000 (can be increased)
- Deployment:
 - Lambda function deployment size (compressed .zip): 50 MB
 - Size of uncompressed deployment (code + dependencies): 250MB
 - Can use the /tmp directory to load other files at startup
 - Size of environment variables: 4 KB

Customization At The Edge

- Many modern applications execute some form of the logic at the edge
- Edge Function:**
 - A code that you write and attach to CloudFront distributions
 - Runs close to your users to minimize latency
- CloudFront provides two types: **CloudFront Functions & Lambda@Edge**
- You don't have to manage any servers, deployed globally
- Use case: customize the CDN content
- Pay only for what you use
- Fully serverless

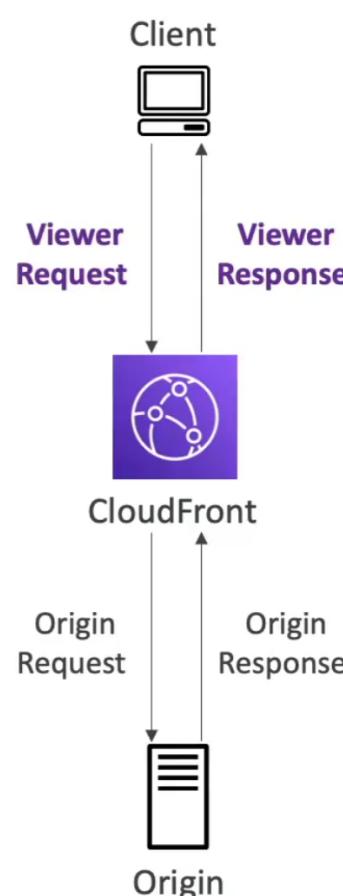
Use Cases

- Website Security and Privacy
- Dynamic Web Application at the Edge
- Search Engine Optimization (SEO)
- Intelligently Route Across Origins and Data Centers
- Bot Mitigation at the Edge

- Real-time Image Transformation
- A/B Testing
- User Authentication and Authorization
- User Prioritization
- User Tracking and Analytics

CloudFront Functions

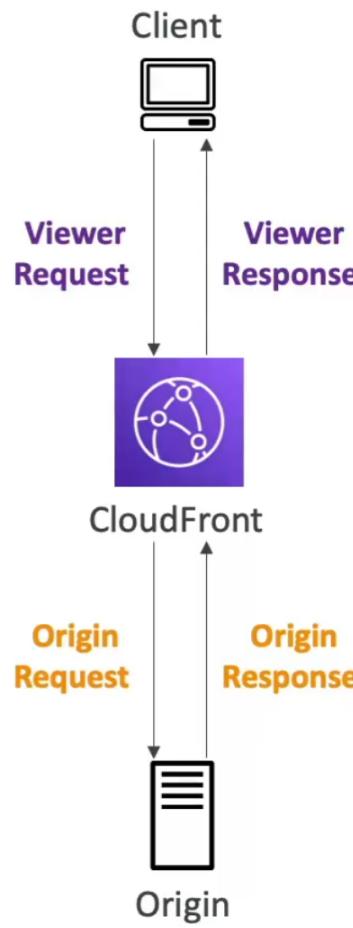
- Lightweight functions written in JavaScript
- For high-scale, latency-sensitive CDN customizations
- Sub-ms startup times, millions of requests/second
- Used to change Viewer requests and responses:
 - **Viewer Request:** after CloudFront receives a request from a viewer
 - **Viewer Response:** before CloudFront forwards the response to the viewer
- Native feature of CloudFront (manage code entirely within CloudFront)



Lambda@Edge

- Lambda functions written in NodeJS or Python
- Scales to 1000s of requests/second
- Used to change CloudFront requests and responses:
 - **Viewer Request** - after CloudFront receives a request from a viewer
 - **Origin Request** - before CloudFront forwards the request to the origin
 - **Origin Response** - after CloudFront receives the response from the origin
 - **Viewer Response** - before CloudFront forwards the response to the viewer

- Author your functions in one AWS Region (us-east-1), then CloudFront replicates to its locations



CloudFront Functions vs. Lambda@Edge - Use Cases

CloudFront Functions

- Cache key normalization
 - Transform request attributes (headers, cookies, query strings, URL) to create an optimal Cache Key
- Header manipulation
 - Insert/modify/delete HTTP headers in the request or response
- URL rewrites or redirects
- Request authentication & authorization
 - Create and validate user-generated tokens (e.g., JWT) to allow/deny requests

Lambda@Edge

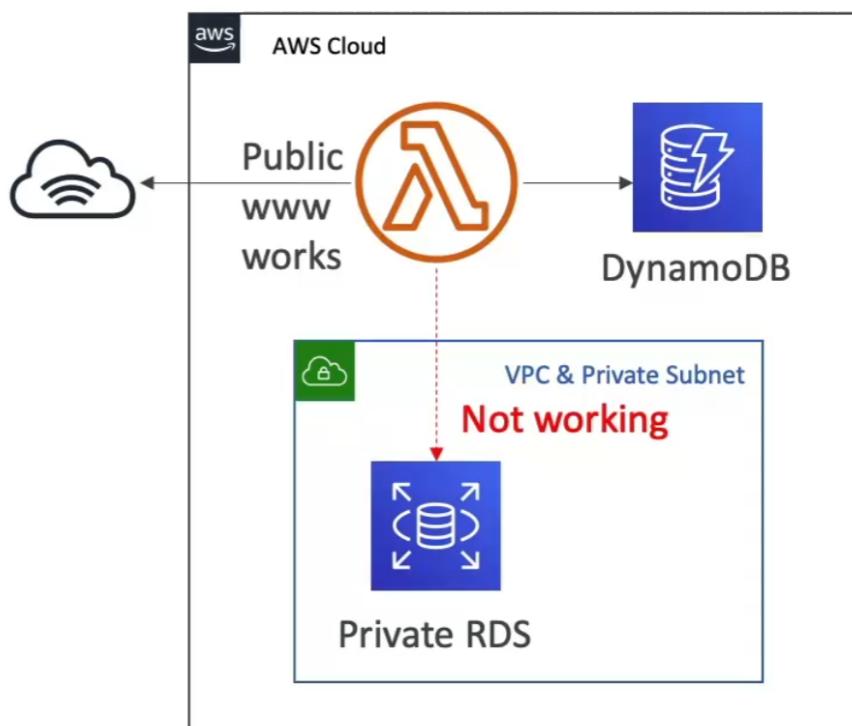
- Longer execution time (several ms)
- Adjustable CPU or memory
- Your code depends on 3rd party libraries (e.g., AWS SDK to access other AWS services)
- Network access to use external services for processing
- File system access or access to the body of HTTP requests

Lambda by default

- By default, your Lambda function is launched outside your own VPC (in an AWS-owned VPC)

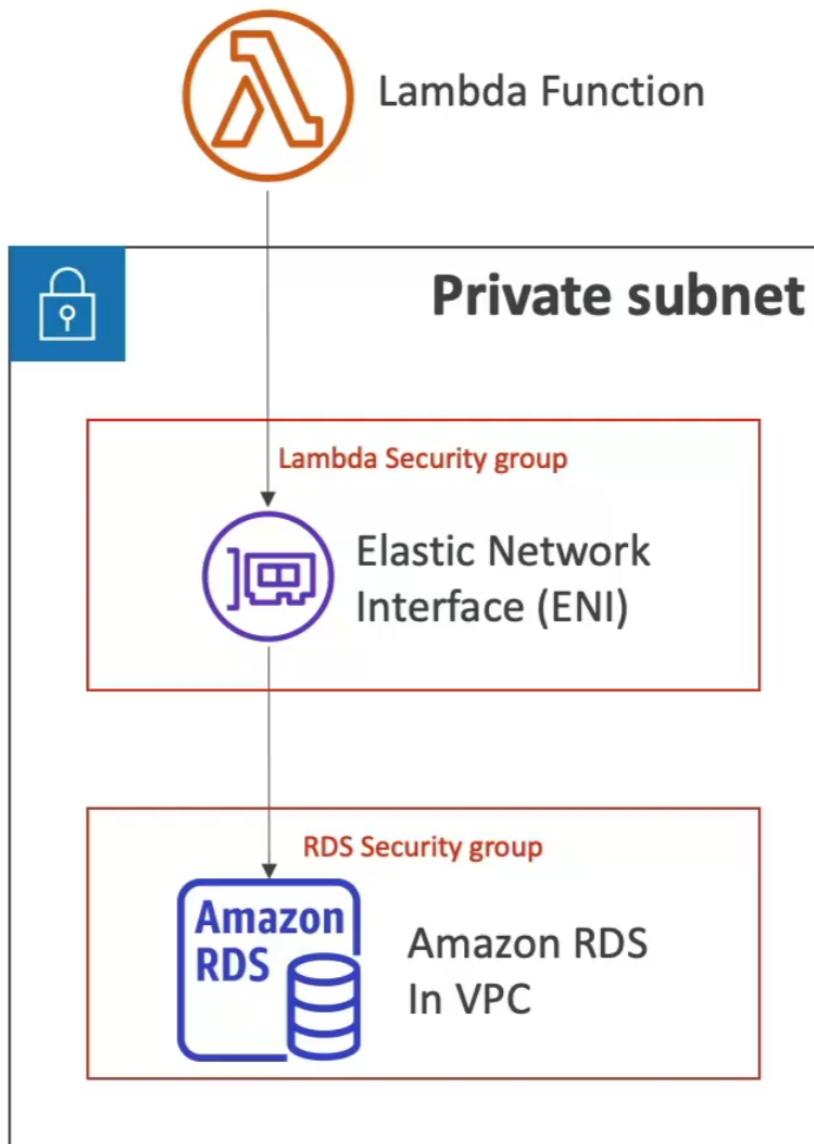
- Therefore, it cannot access resources in your VPC (RDS, ElastiCache, internal ELB...)

Default Lambda Deployment



Lambda in VPC

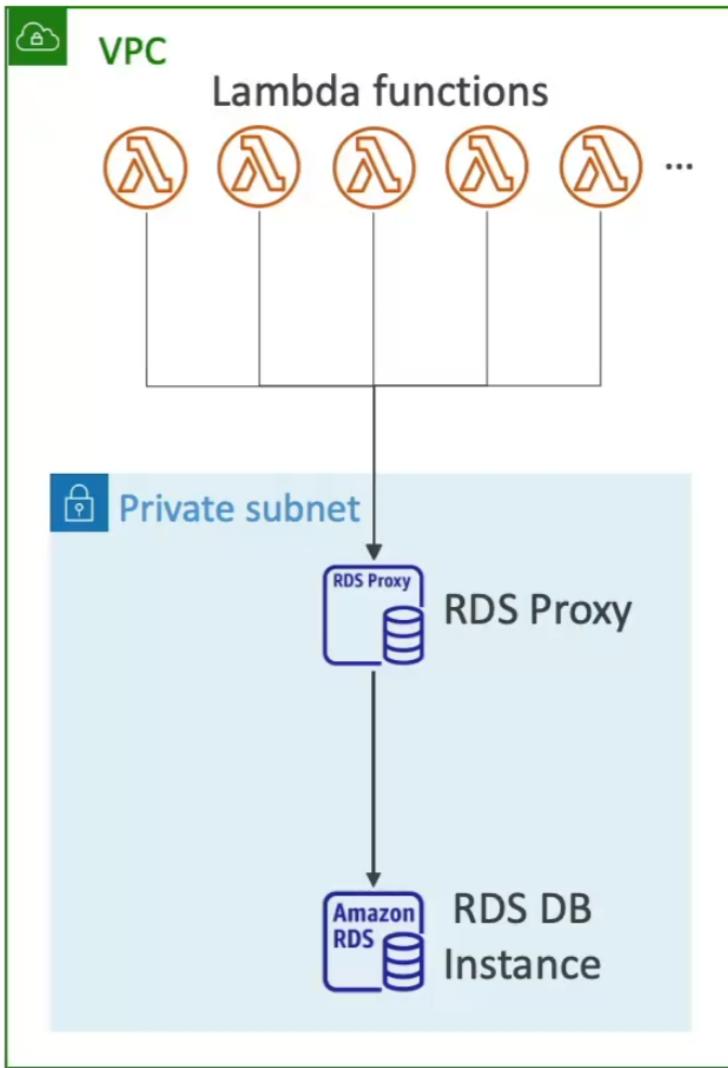
- You must define the VPC ID, the Subnets and the Security Groups
- Lambda will create an ENI (Elastic Network Interface) in your subnets



Lambda with RDS Proxy

- If Lambda functions directly access your database, they may open too many connections under high load
- RDS Proxy
 - Improve scalability by pooling and sharing DB connections
 - Improve availability by reducing by 66% the failover time and preserving connections
 - Improve security by enforcing IAM authentication and storing credentials in Secrets Manager

- The Lambda function must be deployed in your VPC, because RDS Proxy is never publicly accessible



Amazon DynamoDB

- Fully managed, highly available with replication across multiple AZs
- NoSQL database - not a relational database - with transaction support
- Scales to massive workloads, distributed database
- Millions of requests per seconds, trillions of row, 100s of TB of storage
- Fast and consistent in performance (single-digit millisecond)
- Integrated with IAM for security, authorization and administration
- Low cost and auto-scaling capabilities
- No maintenance or patching, always available
- Standard & Infrequent Access (IA) Table Class

DynamoDB - Basics

- DynamoDB is made of **Tables**
- Each table has a **Primary Key** (must be decided at creation time)
- Each table can have an infinite number of items (= rows)
- Each item has **attributes** (can be added over time - can be null)
- Maximum size of an item is **400KB**
- Data types supported are:
 - Scalar Types** - String, Number, Binary, Boolean, Null
 - DocumentTypes** - List, Map
 - Set Types** - String Set, Number Set, Binary Set
- Therefore, in DynamoDB you can rapidly evolve schemas**

Primary Key		Attributes	
Partition Key	Sort Key	Score	Result
User_ID	Game_ID	Score	Result
7791a3d6...	4421	92	Win
873e0634...	1894	14	Lose
873e0634...	4521	77	Win

DynamoDB - Read/Write Capacity Modes

- Control how you manage your table's capacity (read/write throughput)
- **Provisioned Mode (default)**
 - You specify the number of reads/writes per second
 - You need to plan capacity beforehand
 - Pay for provisioned Read Capacity Units (RCU) & Write Capacity Units (WCU)
 - Possibility to add auto-scaling mode for RCU & WCU
- On-Demand Mode
 - Read/write automatically scale up/down with your workloads
 - No capacity planning needed
 - Pay for what you use, more expensive(\$\$\$)
 - Great for unpredictable workloads, steep sudden spikes

Hands on

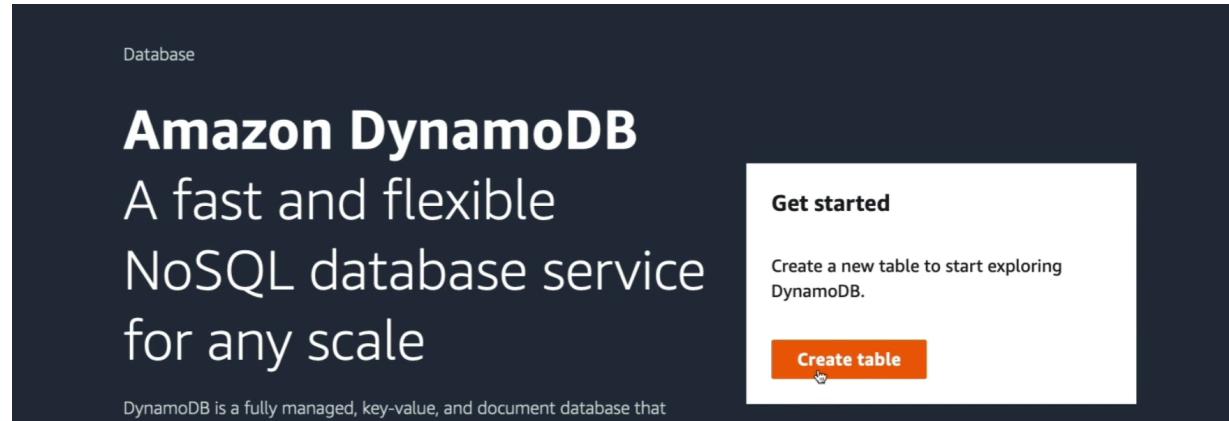


Table details Info

DynamoDB is a schemaless database that requires only a table name and a primary key when you create the table.

Table name
This will be used to identify your table.
 Between 3 and 255 characters, containing only letters, numbers, underscores (_), hyphens (-), and periods (.)

Partition key
The partition key is part of the table's primary key. It is a hash value that is used to retrieve items from your table and allocate data across hosts for scalability and availability.
 1 to 255 characters and case sensitive.

Sort key - optional
You can use a sort key as the second part of a table's primary key. The sort key allows you to sort or search among all items sharing the same partition key.
 1 to 255 characters and case sensitive.

<input type="radio"/> On-demand Simplify billing by paying for the actual reads and writes your application performs.	<input checked="" type="radio"/> Provisioned Manage and optimize your costs by allocating read/write capacity in advance.
---	---

Read capacity

Auto scaling [Info](#)

Dynamically adjusts provisioned throughput capacity on your behalf in response to actual traffic patterns.

- On
- Off

Provisioned capacity units

1

Write capacity

Auto scaling [Info](#)

Dynamically adjusts provisioned throughput capacity on your behalf in response to actual traffic patterns.

- On
- Off

Provisioned capacity units

1

You are not charged any fee for using it.

Encryption key management

Owned by Amazon DynamoDB [Learn more](#)

The key is owned and managed by DynamoDB. You are not charged an additional fee for using this customer master key (CMK).

AWS managed CMK [Learn more](#)

The key is stored in your account and is managed by AWS Key Management Service (AWS KMS). AWS KMS charges apply.

Stored in your account, and owned and managed by you [Learn more](#)

The key is stored in your account and is owned and managed by you. AWS KMS charges apply.

Tags

Tags are pairs of keys and optional values, that you can assign to AWS resources. You can use tags to control access to your resources or track your AWS spending.

No tags are associated with the resource.

[Add new tag](#)

You can add 50 more tags.

Cancel

[Create table](#)

The DemoTable table was created successfully.

DynamoDB > Tables

Tables (1) [Info](#)

Actions ▾ Delete [Create table](#)

Find tables by table name

Any table tag

< 1 > |

<input type="checkbox"/>	Name	Status	Partition key	Sort key	Indexes	Read capacity mode	Write capacity mode
<input type="checkbox"/>	DemoTable		Active	user_id (String)	-	0	Provisioned (1)

DynamoDB > Tables > DemoTable

Partition key	Sort key	Capacity mode	Table status
user_id (String)	-	Provisioned	Active No active alarms
Indexes 0 globals, 0 locals	DynamoDB stream Disabled	Point-in-time recovery (PITR) Disabled	Time to Live (TTL) Info Disabled
Replication Regions 0 Regions	Encryption Owned by Amazon	Date created September 22, 2021, 13:05:23 (UTC+01:00)	

Amazon Resource Name (ARN)

Scan Query

Table or index
DemoTable

Filters

Run Reset

Completed Read capacity units consumed: 0.5

Items returned (0)

Find items

DynamoDB > Items: DemoTable > Item editor

Create item

Form JSON

Attributes

Attribute name	Value	Type
user_id - Partition key	stephane_123	String
name	Stephane Maarek	String
favorite_movie	Memento	String
favorite_number	42	Number

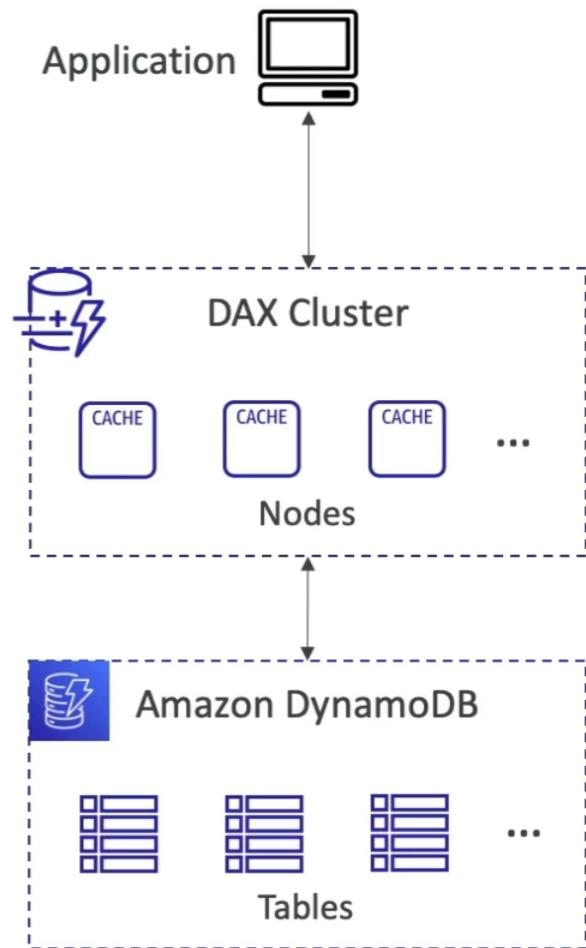
Add new attribute

Cancel Create item

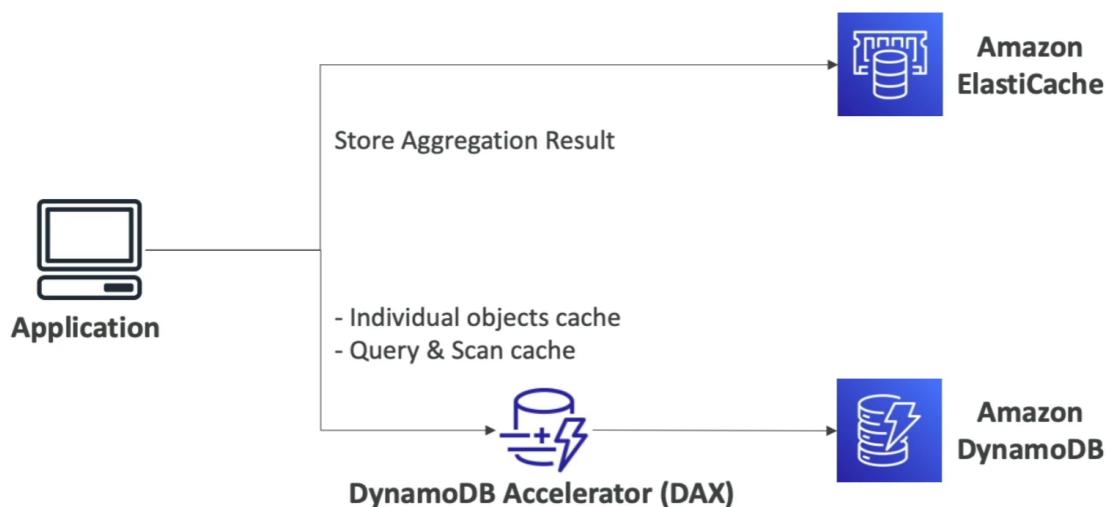
DynamoDB Accelerator (DAX)

- Fully-managed, highly available, seamless in-memory cache for DynamoDB
- **Help solve read congestion by caching**
- **Microseconds latency for cached data**
- Doesn't require application logic modification (compatible with existing DynamoDB APIs)

- 5 minutes TTL for cache (default)



DynamoDB Accelerator (DAX) vs. ElastiCache

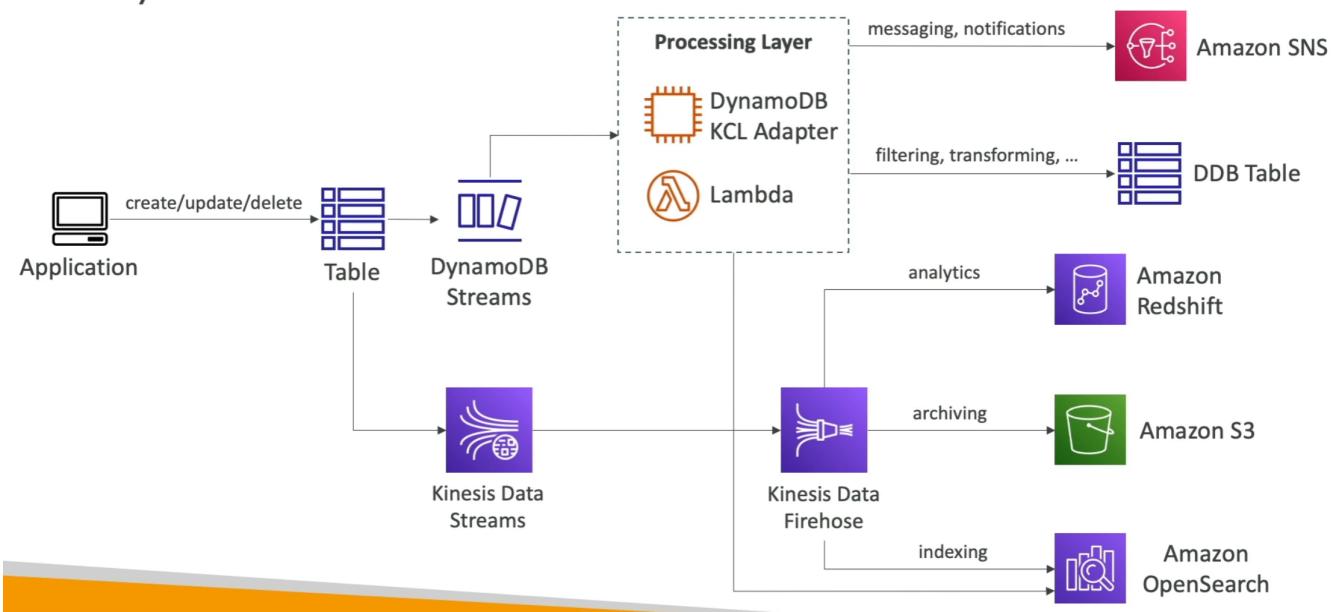


DynamoDB - Stream Processing

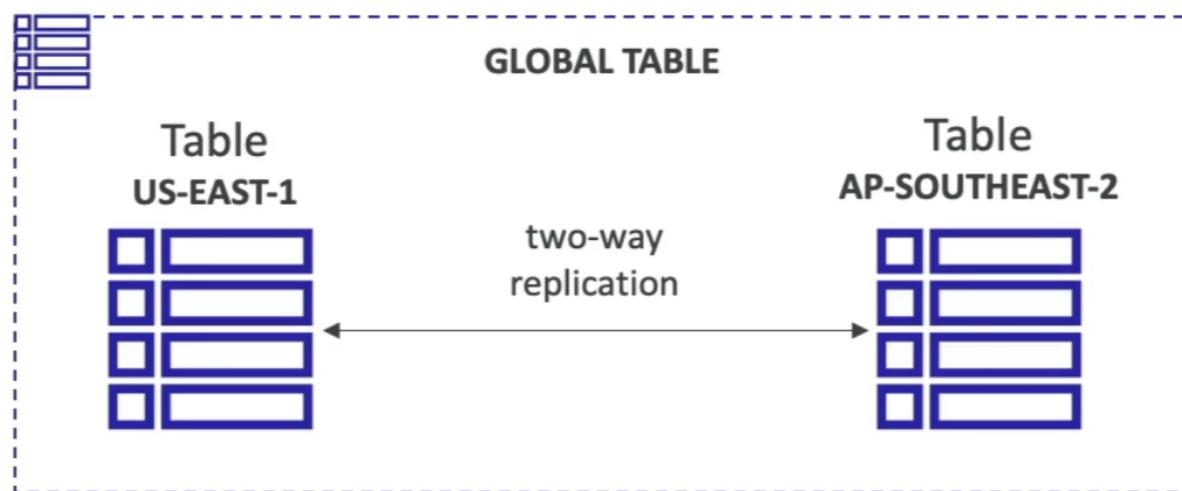
- Ordered stream of item-level modifications (create/update/delete) in a table
- Use cases:
 - React to changes in real-time (welcome email to users)
 - Real-time usage analytics
 - Insert into derivative tables
 - Implement cross-region replication
 - Invoke AWS Lambda on changes to your DynamoDB table

DynamoDB Streams	Kinesis Data Streams (newer)
24 hours retention	1 year retention
Limited # of consumers	High # of consumers
Process using AWS Lambda Triggers, or DynamoDB Stream Kinesis adapter	Process using AWS Lambda, Kinesis Data Analytics, Kinesis Data Firehose, AWS Glue Streaming ETL ...

DynamoDB Streams



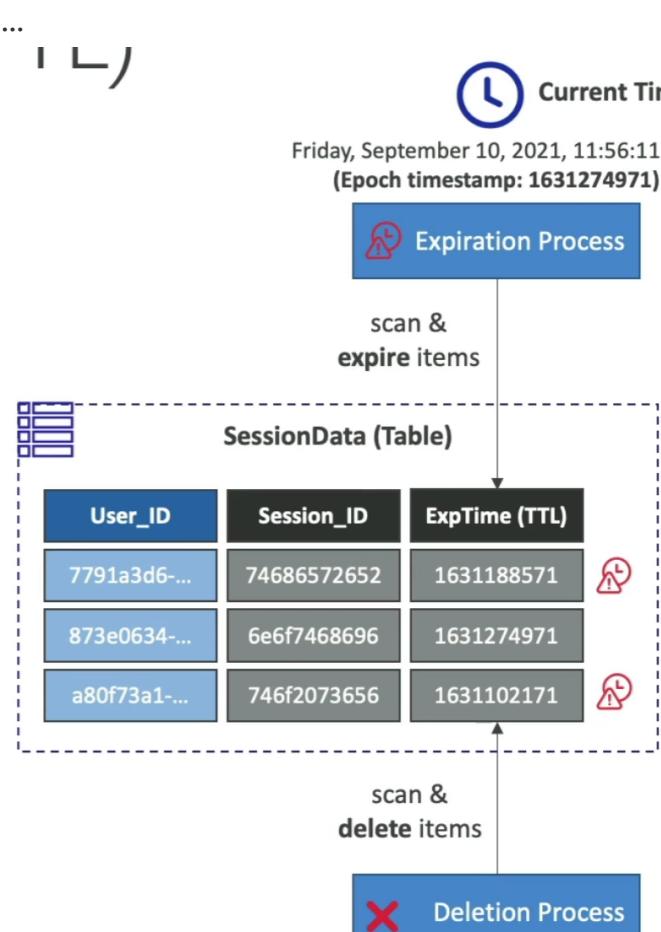
DynamoDB Global Tables



- Make a DynamoDB table accessible with **low latency** in multiple-regions
- Active-Active replication
- Applications can READ and WRITE to the table in any region
- Must enable DynamoDB Streams as a pre-requisite

DynamoDB - Time To Live (TTL)

- Automatically delete items after an expiry timestamp
- Use cases: reduce stored data by keeping only current items, adhere to regulatory obligations, web session handling



DynamoDB - Backups for disaster recovery

- Continuous backups using point-in-time recovery (PITR)
 - Optionally enabled for the last 35 days
 - Point-in-time recovery to any time within the backup window
 - The recovery process creates a new table
- On-demand backups
 - Full backups for long-term retention, until explicitly deleted
 - Doesn't affect performance or latency
 - Can be configured and managed in AWS Backup (enables cross-region copy)
 - The recovery process creates a new table

DynamoDB - Integration with Amazon S3

- Export to S3 (must enable PITR)
 - Works for any point of time in the last 35 days
 - Doesn't affect the read capacity of your table
 - Perform data analysis on top of DynamoDB
 - Retain snapshots for auditing
 - ETL on top of S3 data before importing back into DynamoDB
 - Export in DynamoDB JSON or ION format



- Import from S3
 - Import CSV, DynamoDB JSON or ION format
 - Doesn't consume any write capacity
 - Creates a new table
 - Import errors are logged in CloudWatch Logs

