



# Serverless

The AWS-related services

# Serverless?

- Debunking the myths:

*Serverless does not mean there are no servers*

- It means, the developer:
  - Doesn't need to care about servers when coding.
  - Doesn't have to manage physical capacity.

# Without Serverless ....

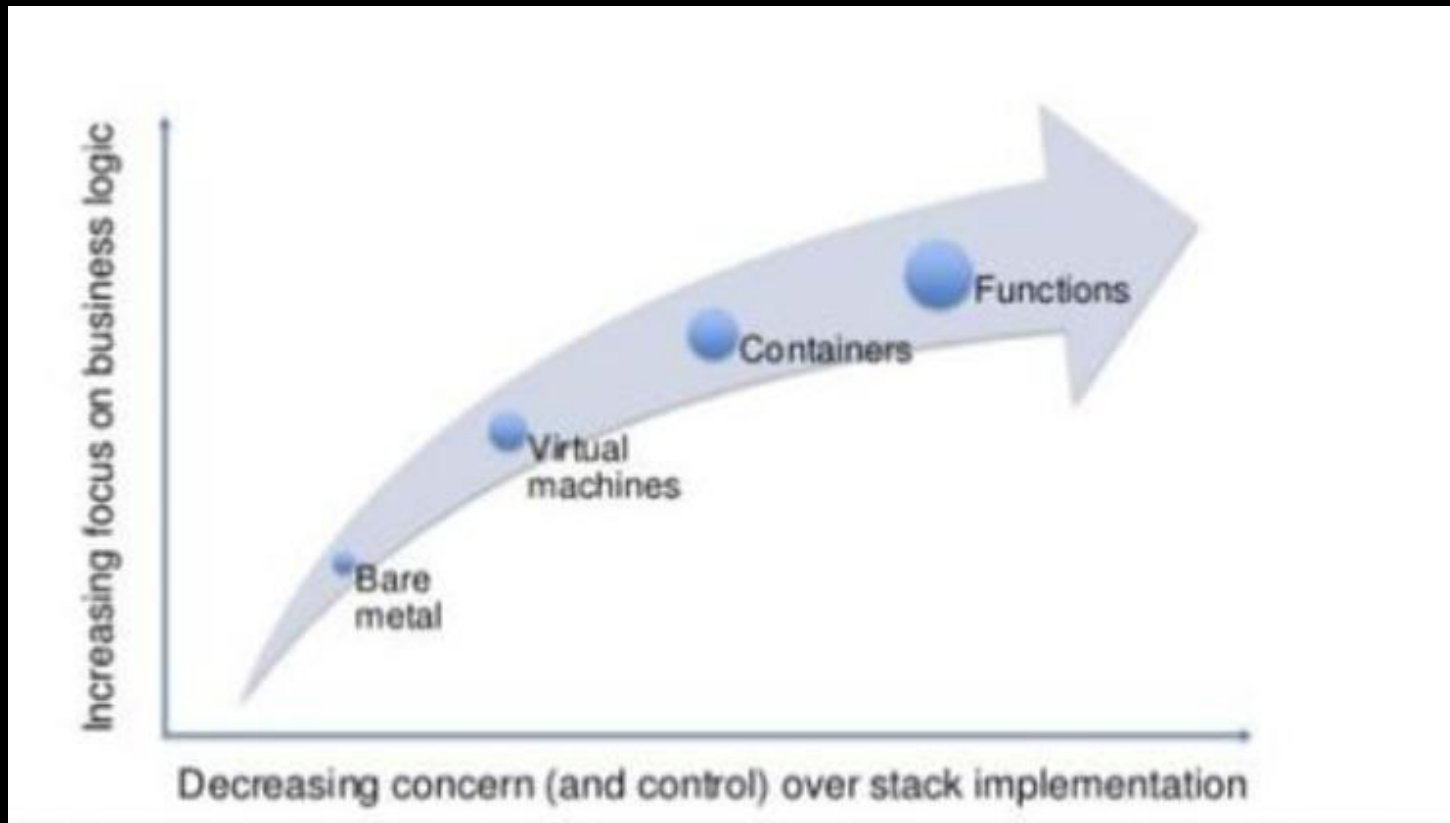
- What is the right size for my server?
  - What capacity is left on my server?
  - How many servers should I provision?
  - How do I handle hardware failures?
- 
- What OS should my server run?
  - Whom should have access to my servers?
  - How do I detect a compromised server?
  - How do I keep my server patched?

# With Serverless

- No servers to provision or manage.
- Scales automatically based on usage.
- Pay for how much you use.
- Availability and fault tolerance built-in.
- Benefits:
  1. Greater agility.
  2. Less overhead.
  3. Increased scale.
  4. Better focus.
  5. Faster time to market.

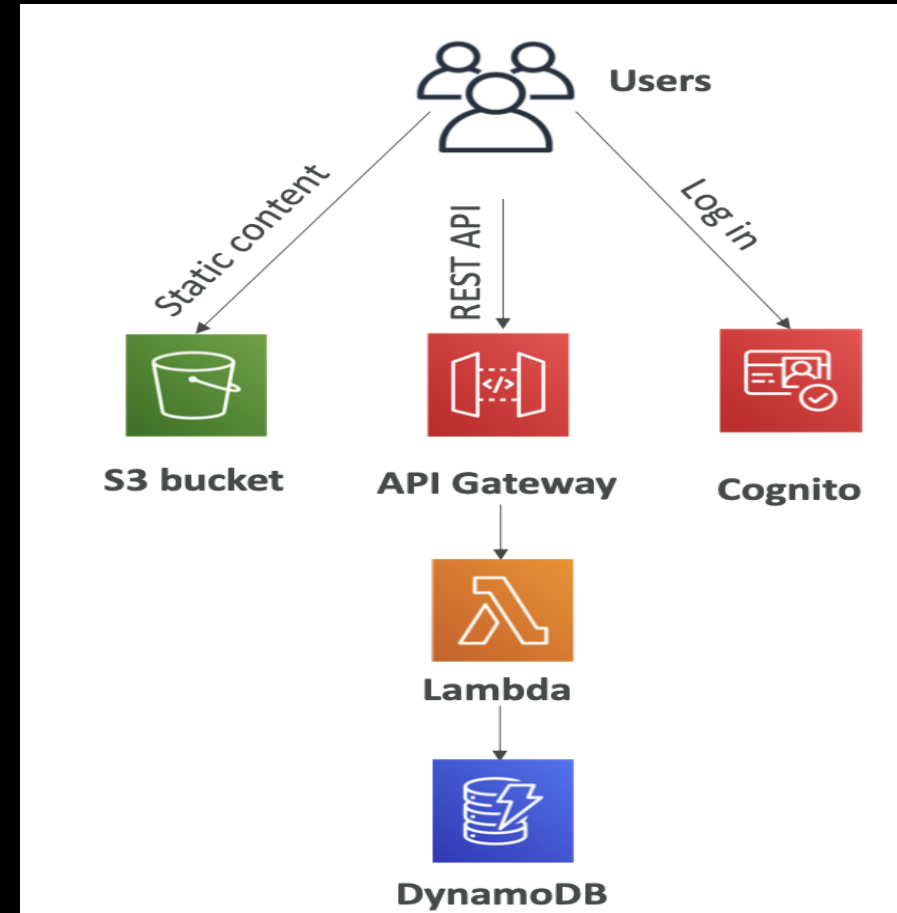
# Developer focus.

- Developers should focus on the product, not infrastructure.



# Serverless Services on AWS

- AWS Lambda. (Compute)
- DynamoDB. (NoSQL)
- AWS Cognito. (User Accounts Mgt.)
- API Gateway (HTTP/REST endpoints)
- S3 (Storage)
- SNS & SQS. (Messaging)
- AWS Kinesis Data Firehose
- Aurora Serverless (RDB)
- Step Functions (Orchestration)
- Fargate (Containers)
- And more ...





# AWS Lambda Service

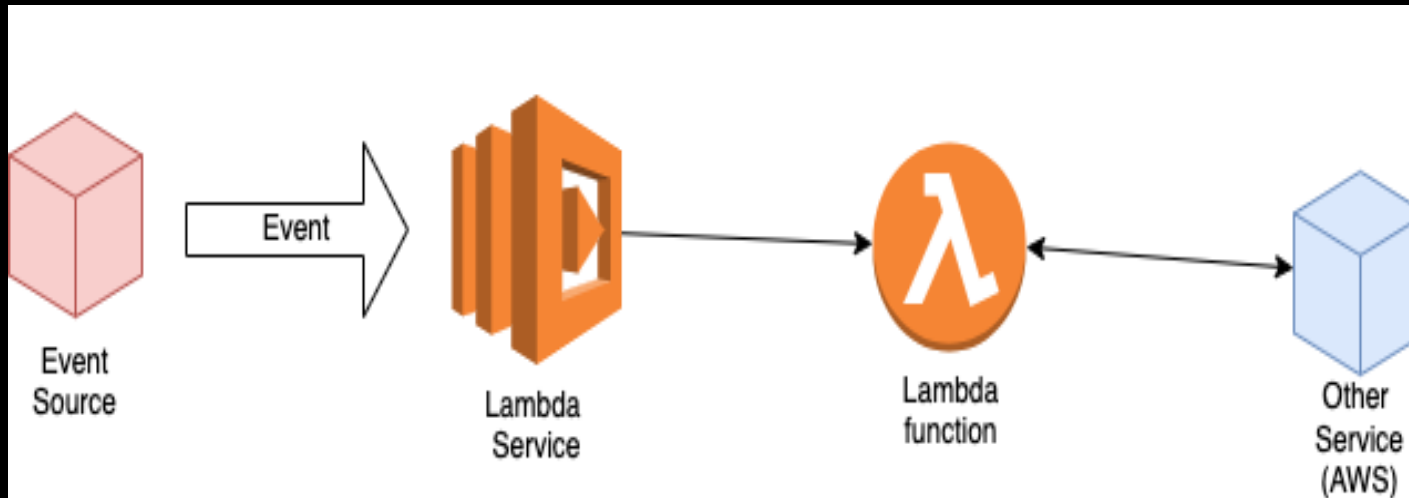
Serverless Compute

# AWS Lambda

- “Lambda is an event-driven, serverless computing platform provided by AWS. It is a computing service that runs code (a function) in response to events and automatically manages the computing resources (CPU, memory, networking) required by that code. It was introduced on November 13, 2014.” Wikipedia
- “Custom code that runs in an ephemeral container” Mike Robins
- FaaS (Functions as a Service)
  - IaaS, PaaS, SaaS



# Lambda runtime model.



- Event Source (Trigger): HTTP request; Change in data store, e.g. database, S3; Change in resource state, e.g. EC2 instance.
- Lambda function: Python, Node, Java, Go, C#, etc

# AWS Lambda service

- The Lambda service manages:
  1. Auto scaling (horizontal)
  2. Load balancing.
  3. OS maintenance.
  4. Security isolation.
  5. Utilization (Memory, CPU, Networking)
- Characteristics:
  1. Function as a unit of scale.
  2. Function as a unit of deployment.
  3. Stateless nature.
  4. Limited by time - short executions.
  5. Run on-demand.
  6. Pay per execution and compute time – generous free tier.
  7. Do not pay for idle time.

# Anatomy of a Lambda function

```
... imports]..  
[ ..... initialization .....  
  .... e.g. d/b connection ....  
export const handler = async (event, context) => {  
  .....  
};  
  
const localFn = (arg) => {  
  .....  
}
```

- Handler() – function to be executed upon invocation.
- Event object – the payload and metadata provided by the event source.
- Context object – access to runtime information.
- Initialization code executes before the handler; Cold starts only.

# Lambda function configuration

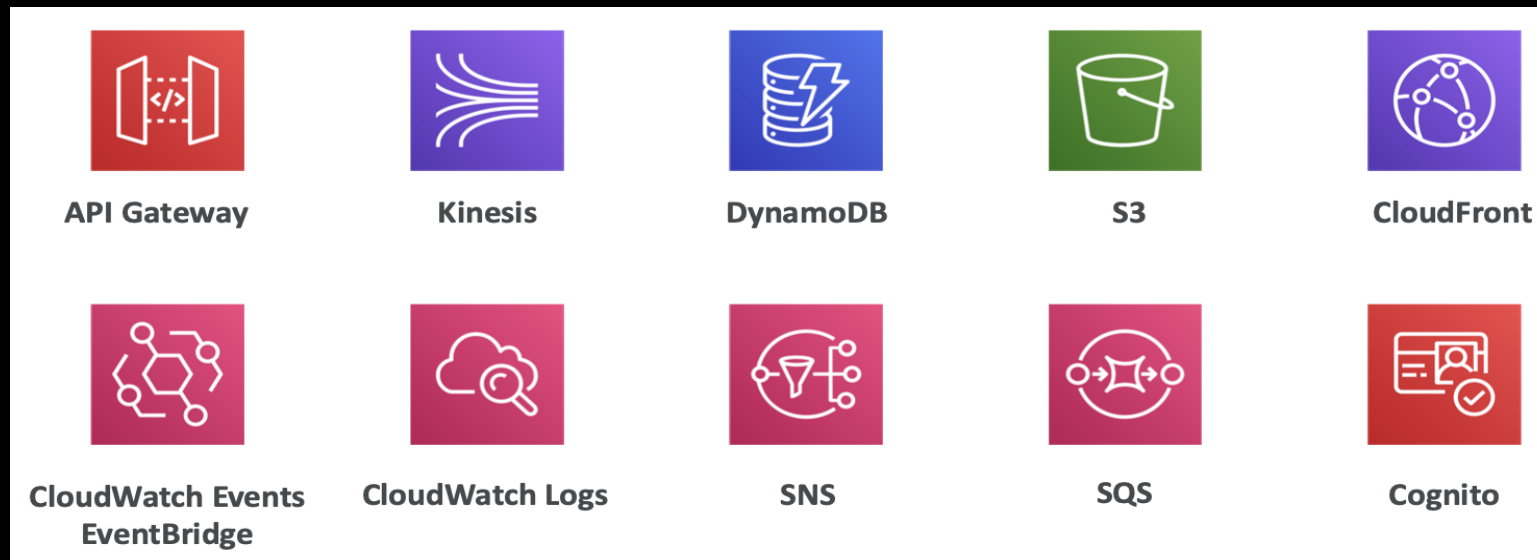
- Lambda service provides a memory control to configure a function's compute power requirements.
  - The % of CPU core and network capacity are computed in proportion to its RAM.
- RAM:
  - From 128MB to 3,008 MB in 64 MB increments
  - The more RAM you add, the more vCPU credits you get.
    - 1,792 MB RAM allocation → one full vCPU reserved.
    - After 1,792 MB → more than one CPU assigned → should use multi-threading to fully utilize.
- For CPU-bound processing, increase the RAM allocation.
- Timeout setting: Max. runtime allowed.
  - Default 3 seconds; maximum 900 seconds (15 minutes).

# Demo

- Objective:
  - Use the CDK to provision a 'Hello World' lambda function.
  - Invoked it (trigger it) using the AWS CLI.
  - See console.log() statement output in the CloudWatch Logs.

# Lambda integration

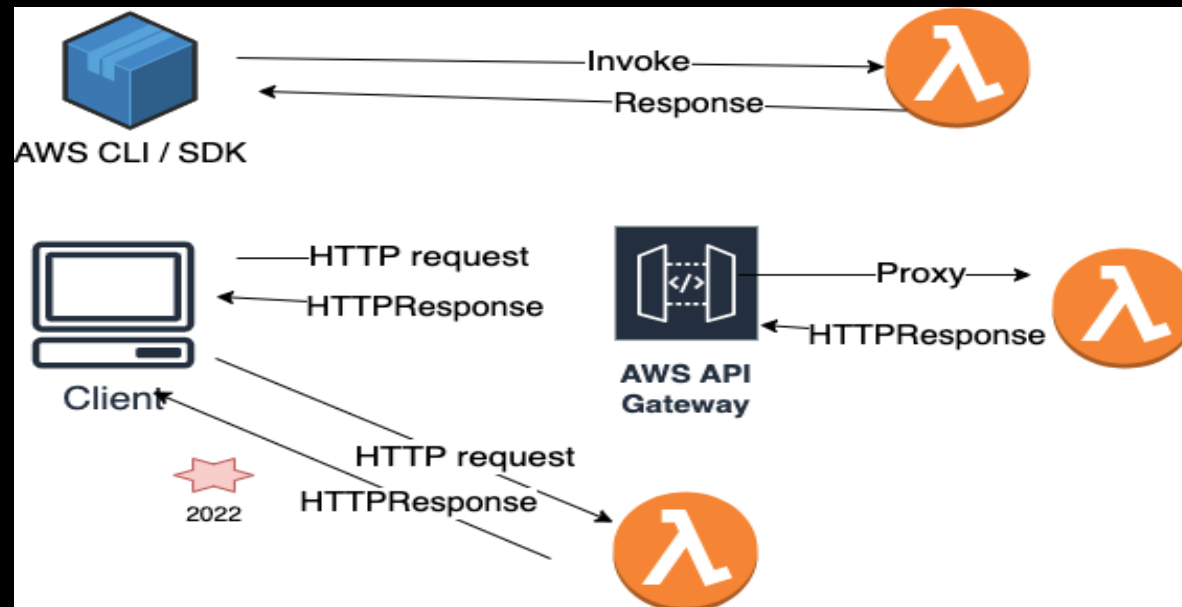
- What services can trigger a lambda function?



- Integration models:
  1. Synchronous.
  2. Asynchronous.
  3. Poll-based

# Synchronous Integration

- Trigger (Event Source): CLI, SDK, API Gateway, Load Balancers, Function URLs.
- Client waits for the response.
- Client should handle error responses (retries, exponential backoff, etc...).



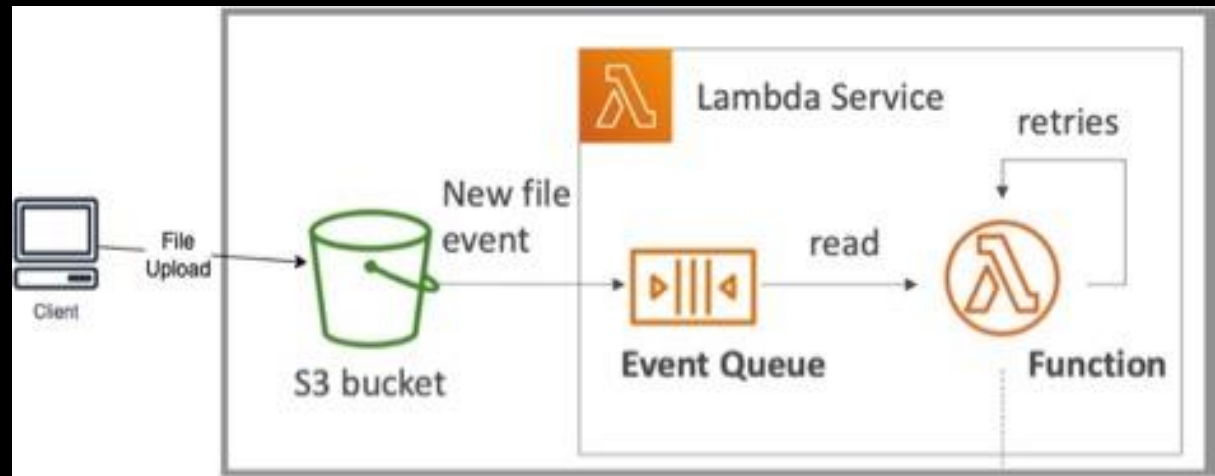
# Demo

- Objective:
  - Use the CDK to:
    1. Provision a lambda function.
    2. Cause Lambda service to generate a URL endpoint.
  - Test URL with Postman client.
  - Configure the endpoint as private and invoke it.



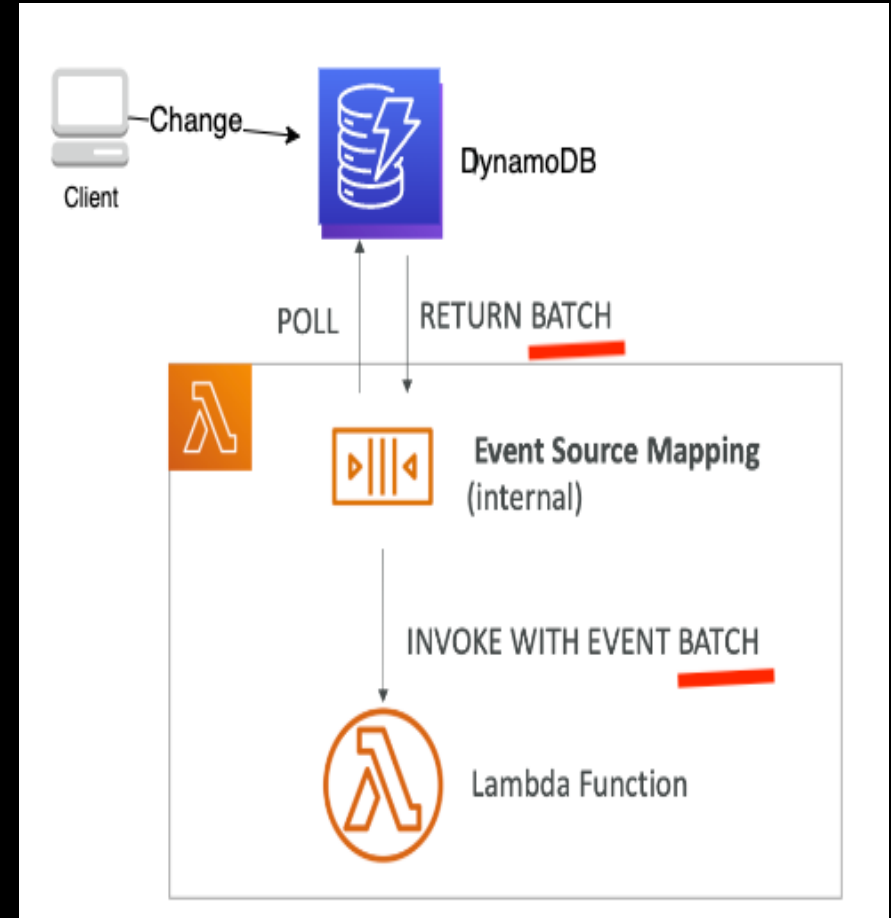
# Asynchronous Integration.

- Trigger (Event Source): S3, SNS, CloudWatch.
- Lambda service places the events in a queue (see below image).
- Lambda service retries on errors – 3 retries, using exponential backoff algorithm.
- Function's processing should be idempotent (due to retries)
- Suitable when application does not require the function result immediately.



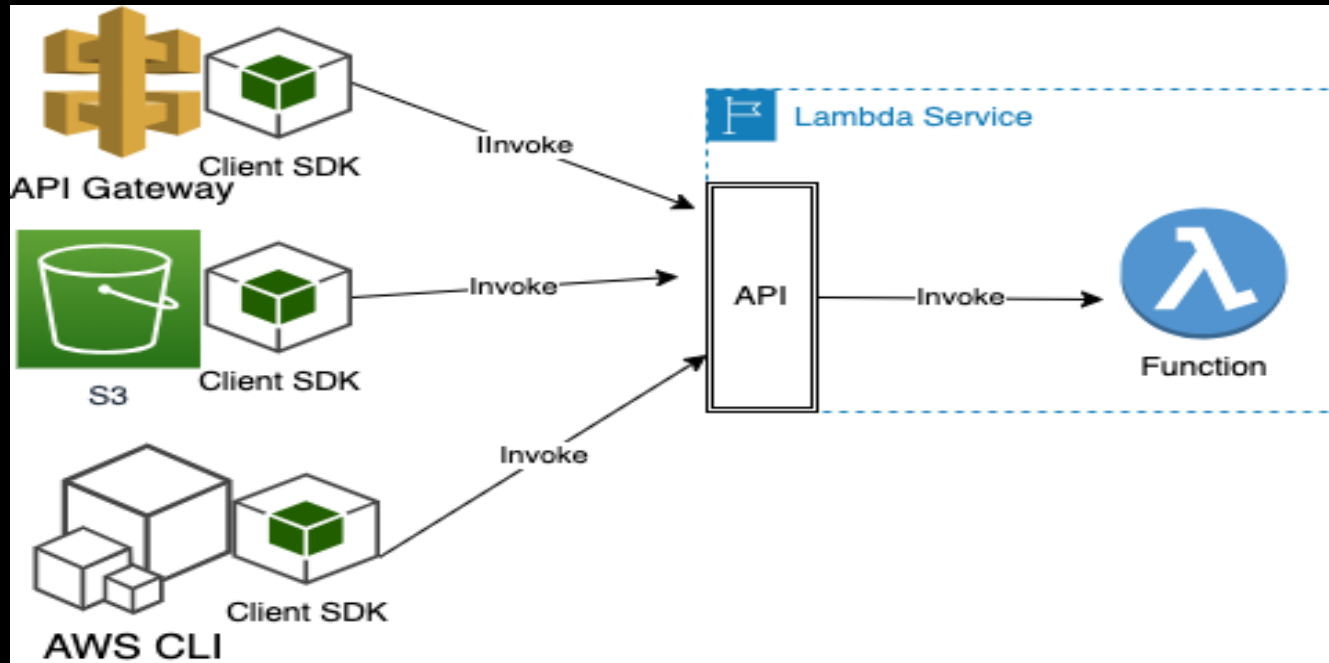
# Event source mapping. (Poll-based Integration)

- Event Sources/Trigger: DynamoDB streams, SQS, Kinesis streams.
- Lambda service polls the source for event records.
- Lambda service invokes the function synchronously.



# Lambda service API & SDK

- Lambda Service provides an API.
- Used by all other services that trigger Lambda functions across all models (sync, async, poll-based).
- Can pass any event payload structure you want.



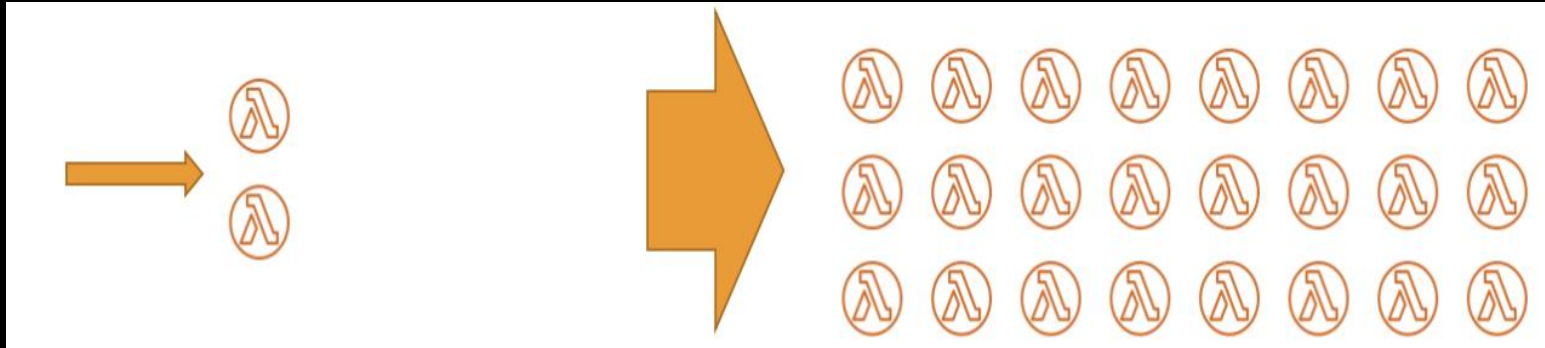
# Execution Role (IAM Role)

- Grants the lambda function permissions to access specified AWS services / resources.
- Many predefined / Managed policies, e.g.
  - AWSLambdaBasicExecutionRole – Upload logs to CloudWatch.
  - AWSLambdaDynamoDBExecutionRole – Read from DynamoDB Streams.
  - AWSLambdaSQSQueueExecutionRole – Read from SQS queue
  - AWSLambdaVPCLambdaAccessExecutionRole – Deploy function in VPC.
- Best practice: Create one Execution Role per function.

# Resource based Policies.

- Use resource-based policies to give other AWS services (and accounts) permission to use your Lambda resource/function.
- An IAM principal (e.g. user, service) can access a Lambda resource:
  - if the IAM policy attached to the principal authorizes it,
  - OR if the function's resource-based policy authorizes it.
- Ex.: An AWS S3 service can trigger a Lambda function if the function's resource-based policy permits it.

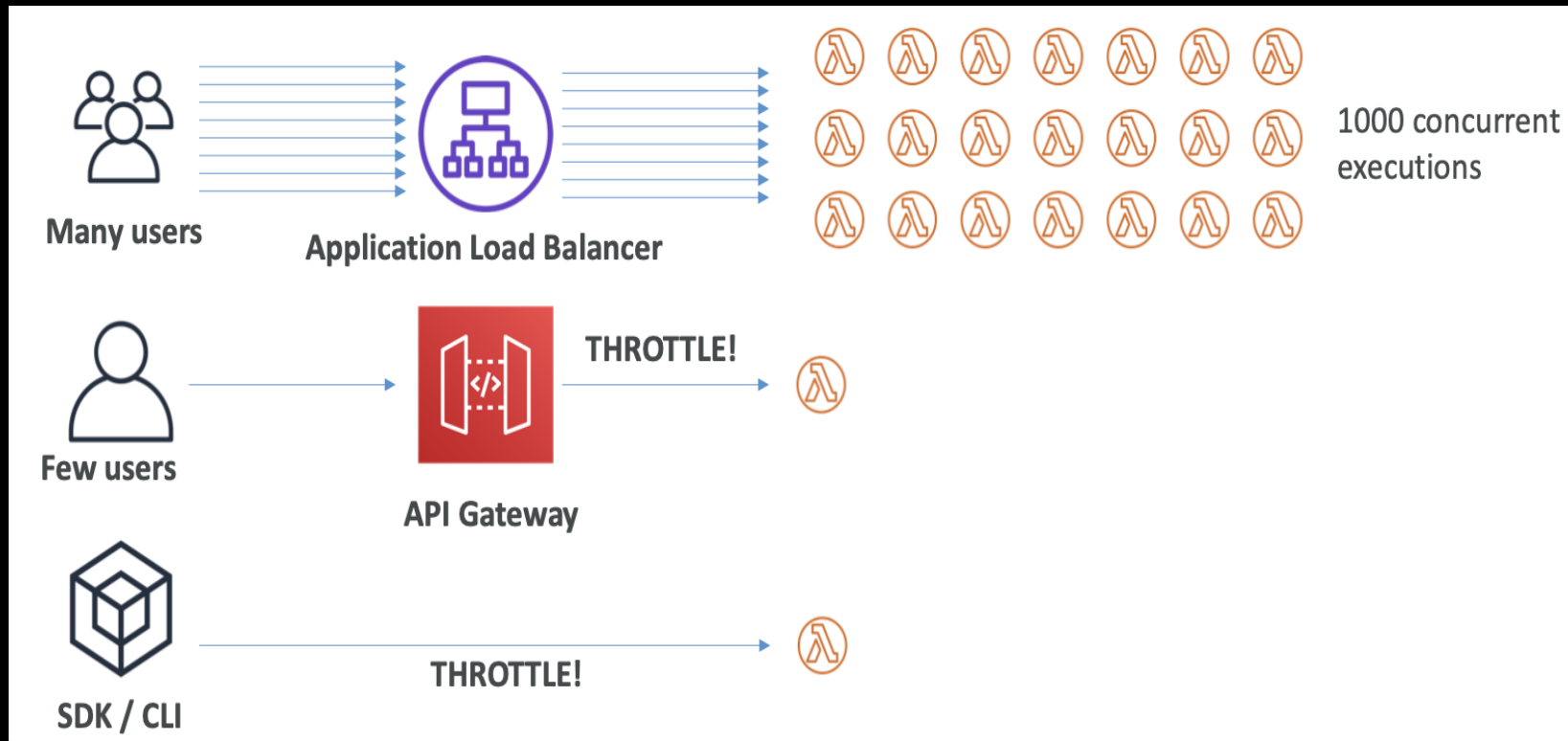
# Concurrency and Throttling



- AWS account concurrency limit is set to 1000 concurrent executions
- Can set a reserved concurrency at the function level (= limit).
  - “Throttle” error response when limit exceeded.
- Throttle behavior:
  - Synchronous invocation => return ThrottleError – HTTP status 429
  - Asynchronous invocation => retry automatically and then go to a DLQ (Dead Letter Queue).

# Concurrency issues

- The following can happen without you reserve (=limit) concurrency:



# Cold Starts & Provisioned Concurrency.

- Cold Start:
  - Source Event → New micro VM instance created → function's initialization code executes → handler code executes
- Warm start – Micro VM reused; Init code not executed.
- Lengthy init code (LoC, dependencies, SDK) effects overall event processing time.
  - Greater latency with cold start executions.
- Solution - Provisioned Concurrency:
  - Micro VMs are pre-allocated in advance; Init code executed during pre-allocation.
  - Cold starts avoided (minimized)
  - lower latency on average.





# AWS DynamoDB Service

Serverless NoSQL Database

# Features

- NoSQL database - not a relational database.
- Fully Managed (Serverless), Highly available with replication across 3 AZ.
- Schema-less (outside the primary key).
  - Records in the same table can have different attributes
- Integrated with IAM for security, authorization and administration.
- Enables event driven programming with DynamoDB Streams

# DynamoDB Basics

- DynamoDB is made up of tables.
- Each table has a primary key (must be decided at creation time).
- A table's entries are termed items (= rows/records).
- Each item has attributes (can be added over time; can be null).
- Maximum size of an item is 400 KB.
- Data types supported are:
  - Scalar Types: String, Number, Binary, Boolean, Null.
  - Document Types: List, Map.
  - Set Types: String Set, Number Set, Binary Set.

# The Primary Key

- Option 1: Simple - Partition key (Hash key) only.
  - Partition key must be unique for each item.
  - Partition key value range must be “diverse” so that the data is distributed evenly.
  - Example: user\_id for a users table.
- Option 2: Composite - Partition key + Sort Key (Range key).
  - The combination must be unique.
  - Item collections - Items are:
    - grouped by partition key.
    - ordered by the sort key.
  - Example: users-games table (Games played by users)
    - user\_id (Partition key) + game\_id (Sort key).

# CDK Table Declaration

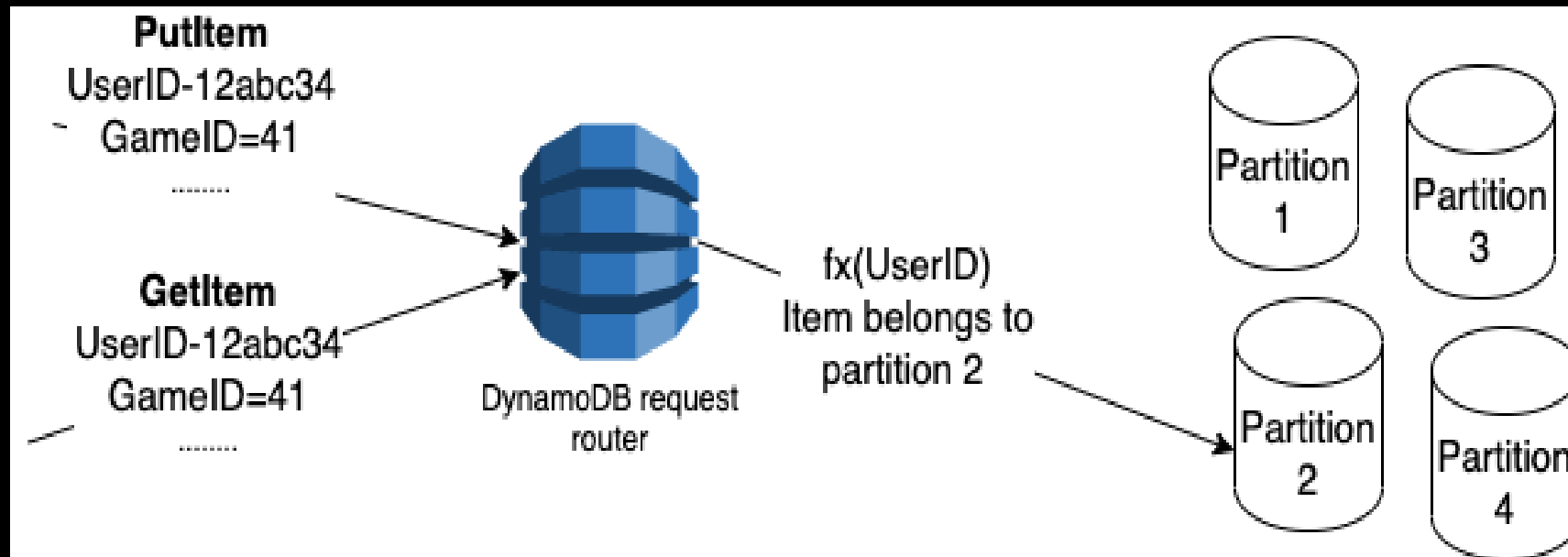
```
7
8 import { RemovalPolicy } from "aws-cdk-lib";
9 import { AttributeType, BillingMode, Table } from "aws-cdk-lib/aws-dynamodb";
10
11 // CDK code for DynamoDB table
12
13 const moviesTable = new Table(this, "MoviesTable", {
14   billingMode: BillingMode.PAY_PER_REQUEST,
15   partitionKey: { name: "movieId", type: AttributeType.NUMBER },
16   removalPolicy: RemovalPolicy.DESTROY,
17   tableName: "Movies",
18 });
19
```

The screenshot displays the AWS Management Console interface for the 'Movies' DynamoDB table. The left-hand navigation pane is open, showing the 'DynamoDB' section with a list of options: Dashboard, Tables, Update settings, Explore items, PartiQL editor, Backups, Exports to S3, Imports from S3, Reserved capacity, and Settings. Red arrows point from the 'Tables' option in the navigation pane to the 'Tables (1)' section in the main content area. Another red arrow points from the 'Explore items' option to the 'Explore table items' button. A third red arrow points from the 'Movies' table entry in the 'Tables (1)' list to the 'General information' section on the right. The 'General information' section shows the following details:

General information	
Partition key	Sort key
movieId (Number)	-
Capacity mode	Table status
On-demand	Active

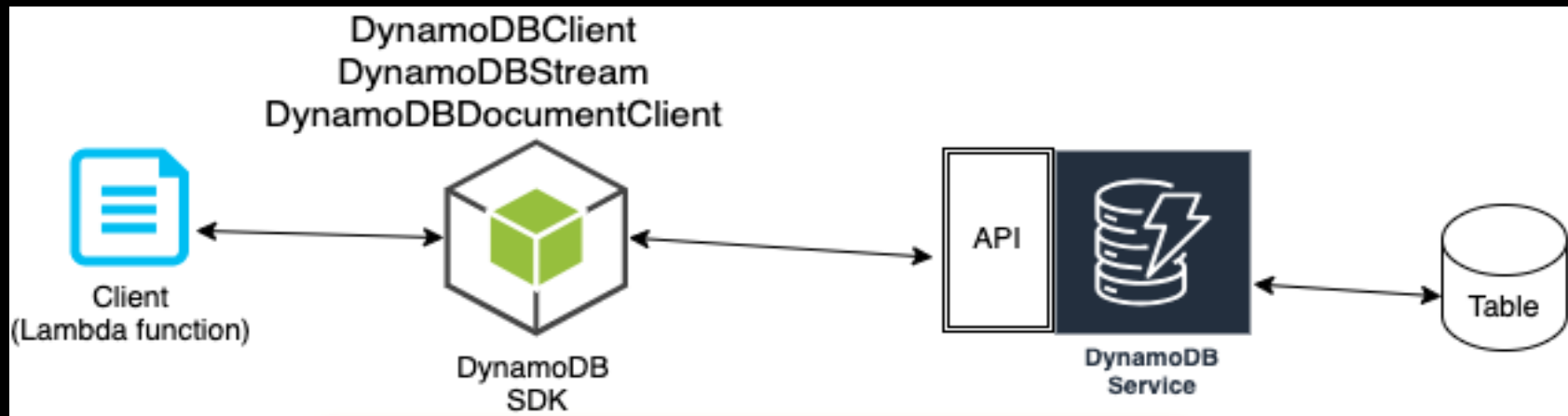
# Horizontal Scaling

- Data is distributed across a fleet of computers, where a single node holds a subset of a table's data, called 'partitions' (max. 10GB).
- Adv – Scalability; Consistent performance.



# DynamoDB API & SDK

- The SDK provides three classes for clients:
  - DynamoDBClient, DynamoDBStreams, and DynamoDBDocumentClient.



```
"dependencies": {  
  "@aws-sdk/client-dynamodb": "^3.67.0",  
  "@aws-sdk/lib-dynamodb": "^3.79.0",  
  "@aws-sdk/util-dynamodb": "^3.303.0",  
  "aws-cdk-lib": "2.71.0",  
  "You, now • Unc
```

# DynamoDB API

- Three main types of actions:
  1. Single-item requests - acts on a single, specific table item and requires the full primary key.
    - PutItem, GetItem, UpdateItem, DeleteItem.
  2. Query - reads a range of items; request must include the partition key.
    - Only suitable for tables with composite primary key.
    - Query result are always from the same partition (item collection).
  3. Scan - reads a range of items but searches across the entire table; an inefficient operation.



# Sample Code

```
7 import { DynamoDBClient } from "@aws-sdk/client-dynamodb";
8 import { DynamoDBDocumentClient, GetCommand } from "@aws-sdk/lib-dynamodb";
9
10 // Native DDB client
11 const ddbClient = new DynamoDBClient({ region: process.env.REGION });
12 // Abstracted DDB client (Document client)
13 const ddbDocClient = DynamoDBDocumentClient.from(
14   ddbClient,
15   ... marshall/unmarshalling options ....
16 );
17
18 const commandOutput = await ddbDocClient.send(
19   new GetCommand({
20     TableName: process.env.TABLE_NAME,
21     Key: { movieId: 1234 },
22   })
23 );
24
```

- **Marshalling and Unmarshalling Data** - The document client allows the use of JS types instead of DynamoDB's native AttributeValues. JS objects passed as parameters are marshalled into AttributeValue shapes required by DDB. Responses from DDB are unmarshalled into plain JS objects.

# API Query actions

- Query requests return items based on:
  - Partition Key value (equals (=) operator only) +  
Sort Key (=, <, <=, >, >=, Between, Begin) – optional.
- Can include a Filter Expression for further filtering (performed on the client side).
- Query response:
  - Up to 1MB of data, or
  - the number of items specified in a Limit parameter.
- Pagination on the response is supported.

- Sample table

Query with these		Not with these	
Primary key		Attributes	
Partition key: CustomerId	Sort key: OrderId	OrderDate	TotalPrice
de91538a	f12801b7	2021-02-12 14:27:21	114.82

# SDK Query Example

- A table stores the cast for each movie, with one item per movie role.

```
24
25 // CDK code for DynamoDB table
26 const movieCastsTable = new Table(this, "MovieCastTable", {
27     billingMode: BillingMode.PAY_PER_REQUEST,
28     partitionKey: { name: "movieId", type: AttributeType.NUMBER },
29     sortKey: { name: "actorName", type: AttributeType.STRING },
30     removalPolicy: RemovalPolicy.DESTROY,
31     tableName: "MovieCast",
32 });
33 //=====
34 // SDK code for DynamoDB query
35 // Find actors whose name begins with Bob on the movie with ID 1234
36 const commandOutput = await ddbDocClient.send(
37     new QueryCommand({
38         TableName: process.env.TABLE_NAME,
39         KeyConditionExpression: "movieId = :m and begins_with(actorName, :a) ",
40         ExpressionAttributeValues: {
41             ":m": 1234,
42             ":a": 'Bob',
43         },
44     })
45 )
```

# Local Secondary Index (LSI)

- LSIs are based on alternate sort key for a table.
- An LSI is stored local to the partition key.
- Max of five LSIs per table.
- Used with Query actions.
- The attribute chosen for LSI sort key must be a scalar String, Number, or Binary.
- LSI must be defined at table creation time.
- Each LSI entry is a projection of the table item.
- A table has:
  - a primary / main index + (optionl) multiple LSIs +  
(optional) multiple GSIs (Global secondary indices)

# Table Indices - Example

- Example: Table of Discussion threads for a social media app.
  - The main / primary index

Partition key		Thread		
ForumName	Sort key Subject	LastPostDateTime	Replies	
"S3"	"aaa"	"2015-03-15:17:24:31"	12	...
"S3"	"bbb"	"2015-01-22:23:18:01"	3	...
"S3"	"ccc"	"2015-02-31:13:14:21"	4	...
"S3"	"ddd"	"2015-01-03:09:21:11"	9	...
"EC2"	"yyy"	"2015-02-12:11:07:56"	18	...
"EC2"	"zzz"	"2015-01-18:07:33:42"	0	...
"RDS"	"rrr"	"2015-01-19:01:13:24"	3	...
"RDS"	"sss"	"2015-03-11:06:53:00"	11	...
"RDS"	"ttt"	"2015-10-22:12:19:44"	5	...
...	...	...	...	

# Table Indices - Example

- Example: Table of Discussion threads for a social media app.
  - A local secondary index (LSI)

<i>LastPostIndex</i>		
<i>Partition key</i>	<i>Sortkey</i>	
<i>ForumName</i>	<i>LastPostDateTime</i>	<i>Subject</i>
"S3"	"2015-01-03:09:21:11"	"ddd"
"S3"	"2015-01-22:23:18:01"	"bbb"
"S3"	"2015-02-31:13:14:21"	"ccc"
"S3"	"2015-03-15:17:24:31"	"aaa"
"EC2"	"2015-01-18:07:33:42"	"zzz"
"EC2"	"2015-02-12:11:07:56"	"yyy"
"RDS"	"2015-01-19:01:13:24"	"rrr"
"RDS"	"2015-02-22:12:19:44"	"ttt"
"RDS"	"2015-03-11:06:53:00"	"sss"
...	...	...

# Information

- Watch <https://www.youtube.com/watch?v=ErPrf74RUDY>