

Serverless Solutions on **AWS**



Nguyen Pham Luan Tien – NashTech Vietnam

May 2023

For NashTech Viet Nam internal circulation only

**Nash
Tech.**

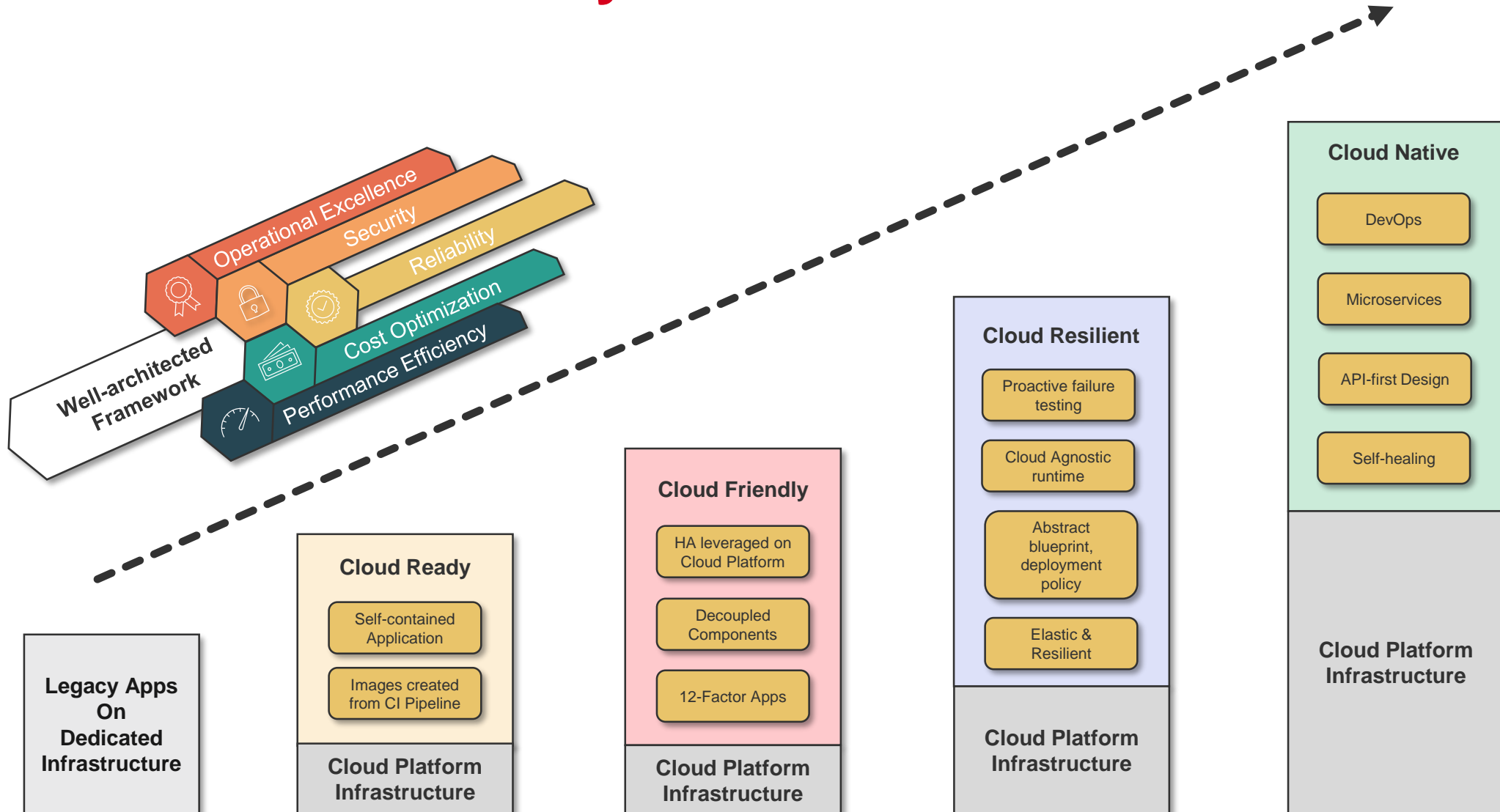
Agenda

1. Why Serverless
2. What is Serverless
3. How to adapt Serverless in AWS
4. Pattern Samples

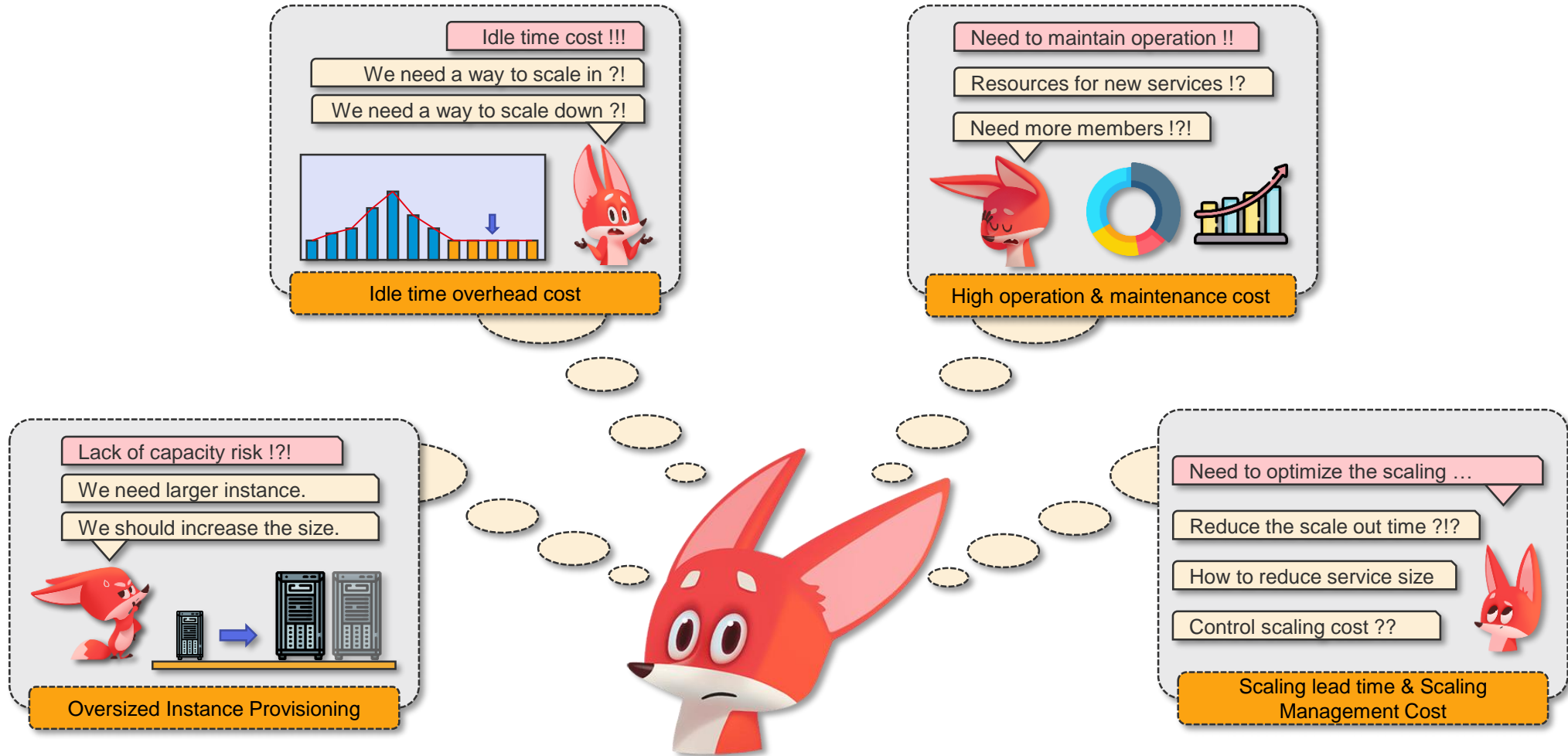
Why Serverless



Cloud-Native Maturity Model



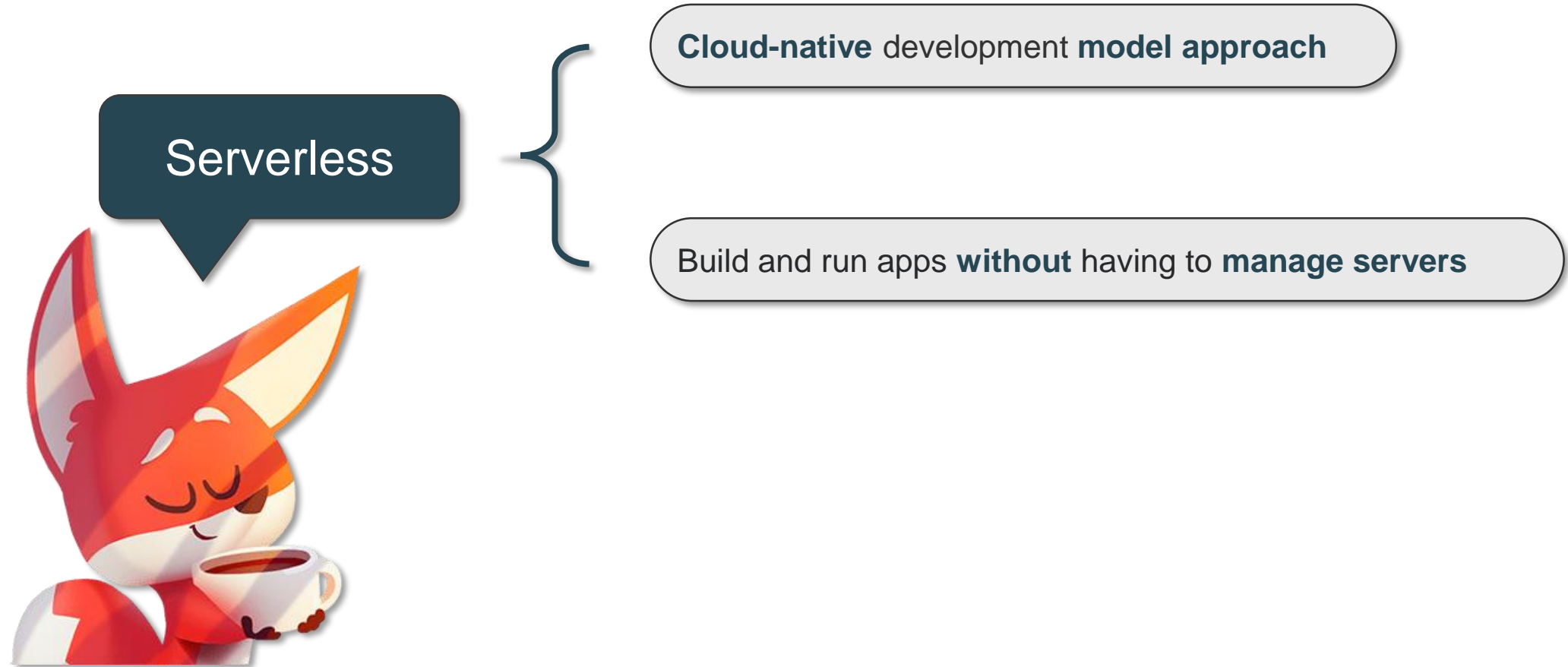
Server hosting Challenges



What is Serverless



Serverless is ...



Serverless has ...

Advantages



Scales on demand **automatically**

No more **Idle time cost**

Eliminates server maintenance and **operation**

Inherently **scalable** with **small** a **light-weight services**

Reduced packaging and **deployment complexity**

Shorten delivery life-cycle, focus right on business

There are ...

Considerations



Vendor lock-in, depends on cloud provider ecosystem

Dealing with cold starts

Testing and debugging become more challenging

Not built for long-running processes

Source code management strategy

Less server control capability

How to adapt Serverless in AWS



Popular serverless services

Popular



CloudFront



API Gateway



Lambda Function



DynamoDB



Amazon Aurora



Amazon Cognito



Simple Storage Service

Integration



Simple Notification Service



Simple Queue Service



Step Function

Monitoring & Analytic



CloudWatch



AWS X-Ray



Kinesis



Athena

Others

Containerization



Elastic Kubernetes Service



Elastic Container Service



Fargate



CodeBuild



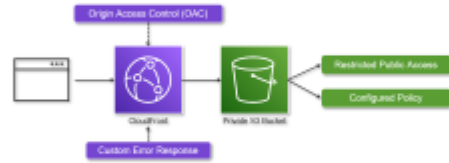
CodeDeploy



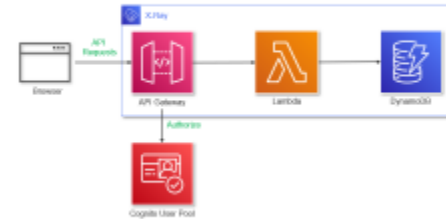
CodePipeline

Common Patterns

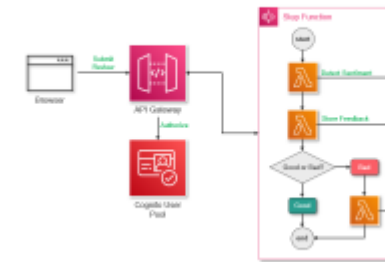
Static Web Hosting



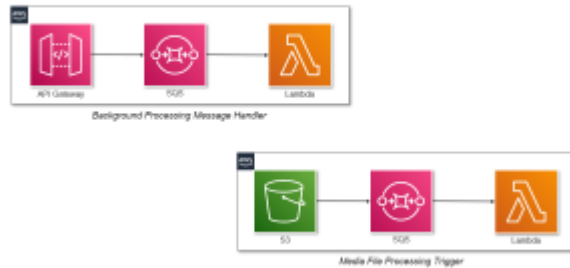
API Service with Lambda



Workflow Service



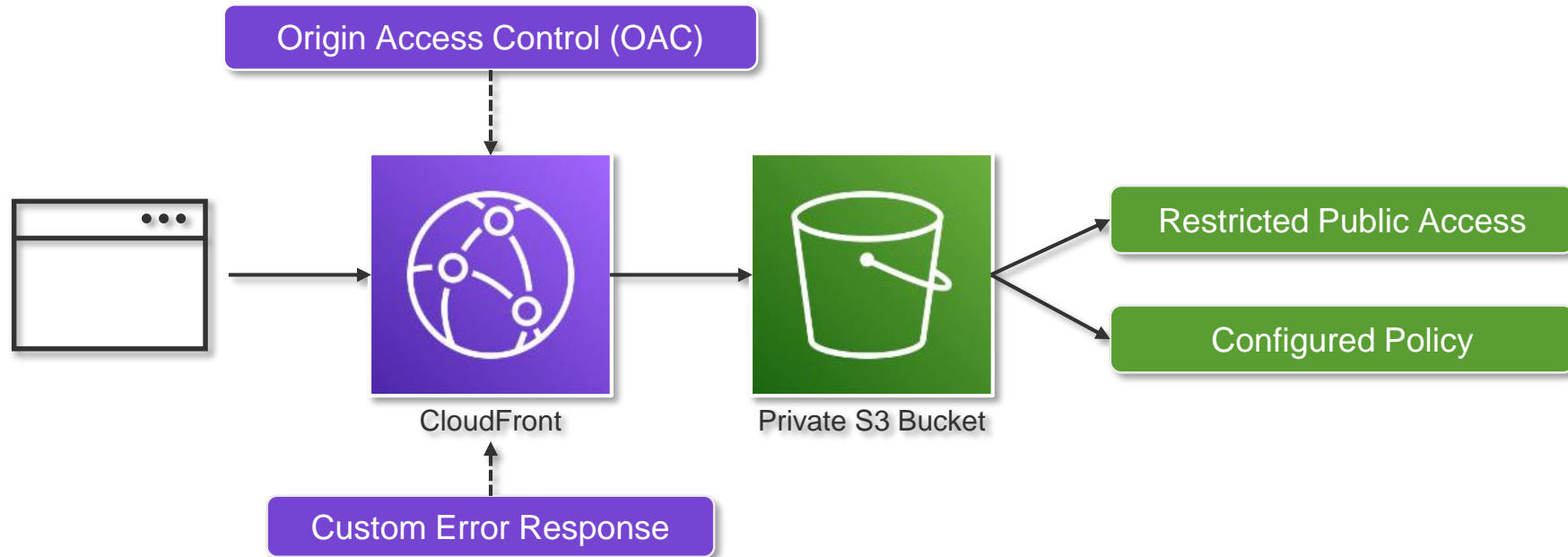
Messaging with SQS



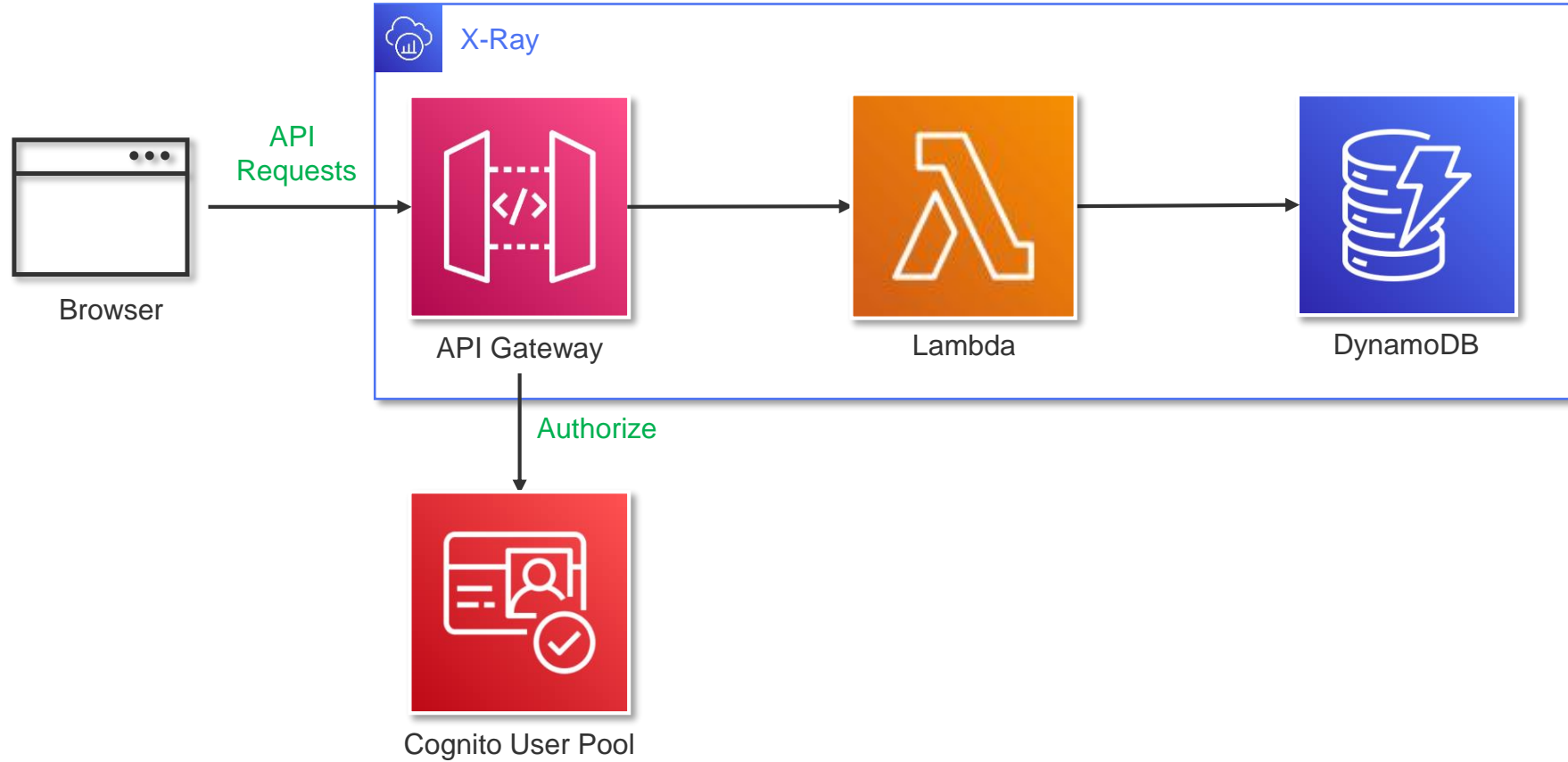
Containerized API Service



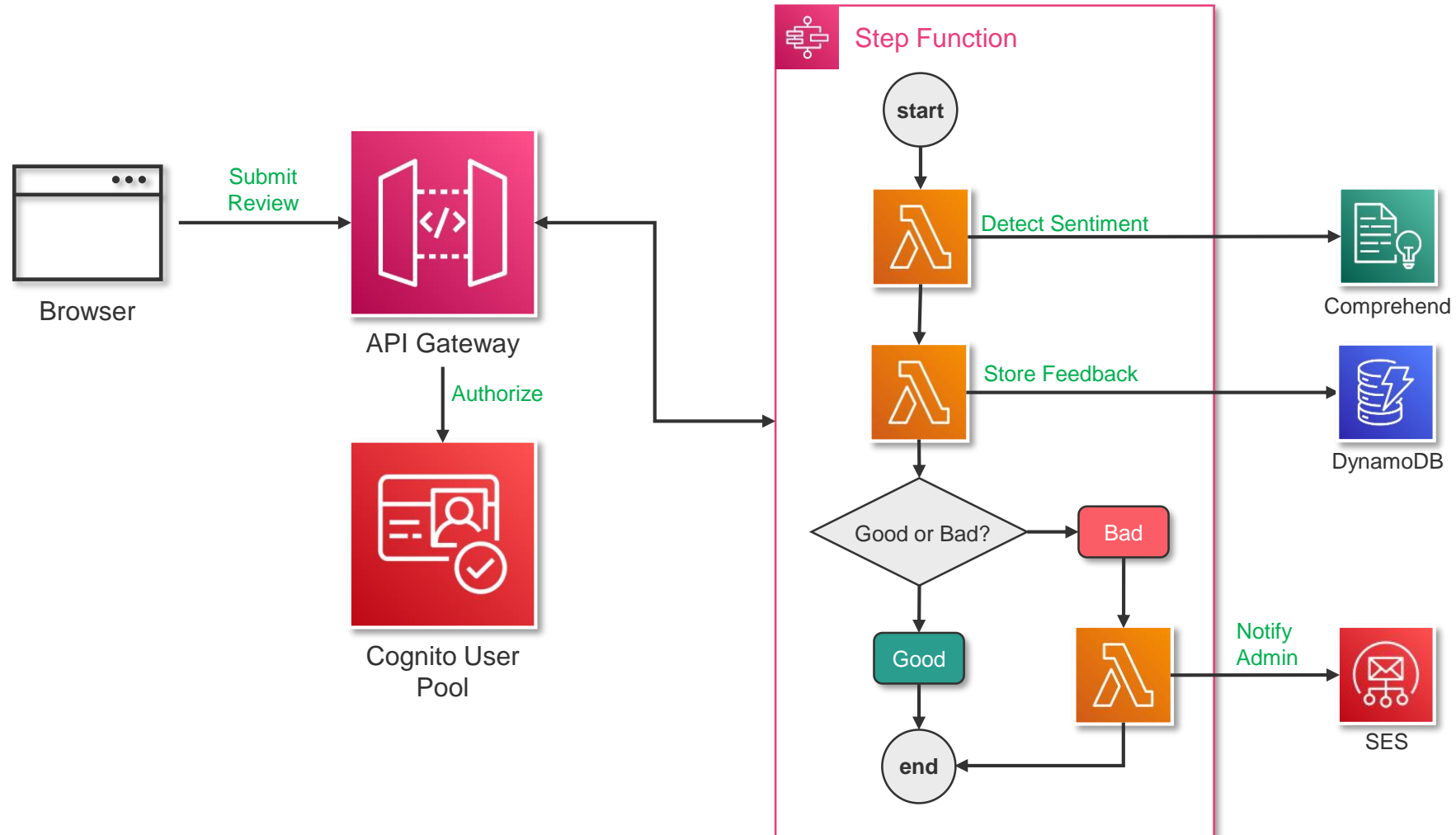
Static Web Hosting



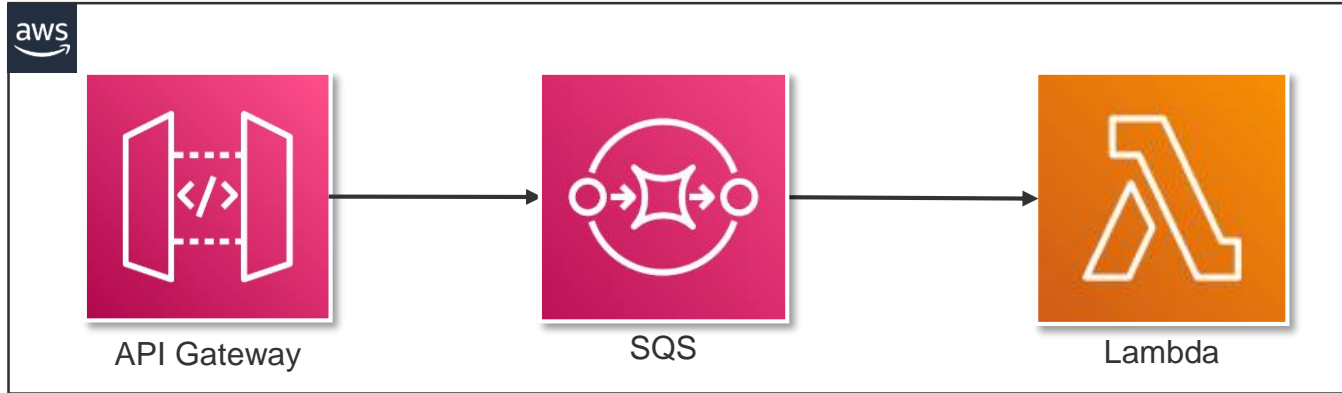
API Service with Lambda



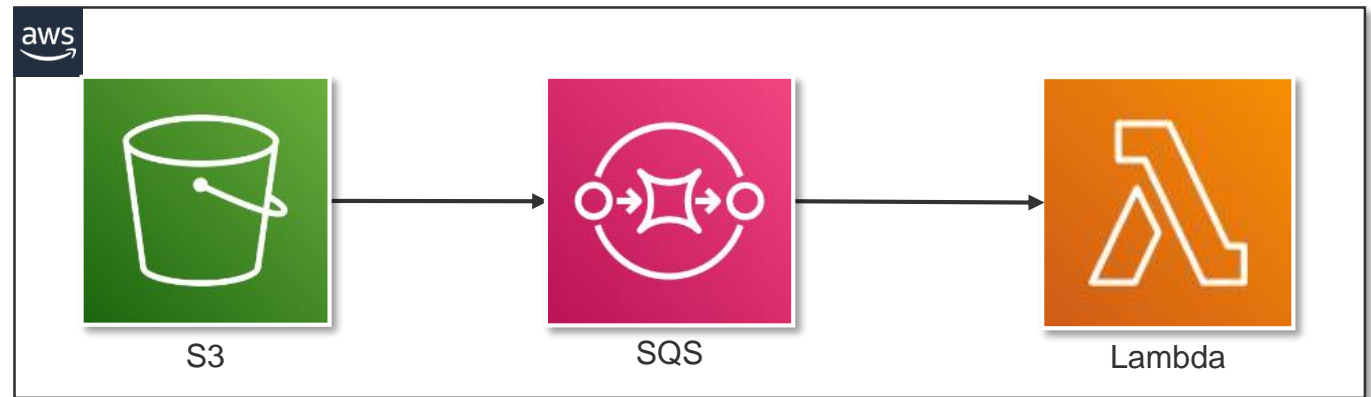
Workflow Service



Messaging with SQS

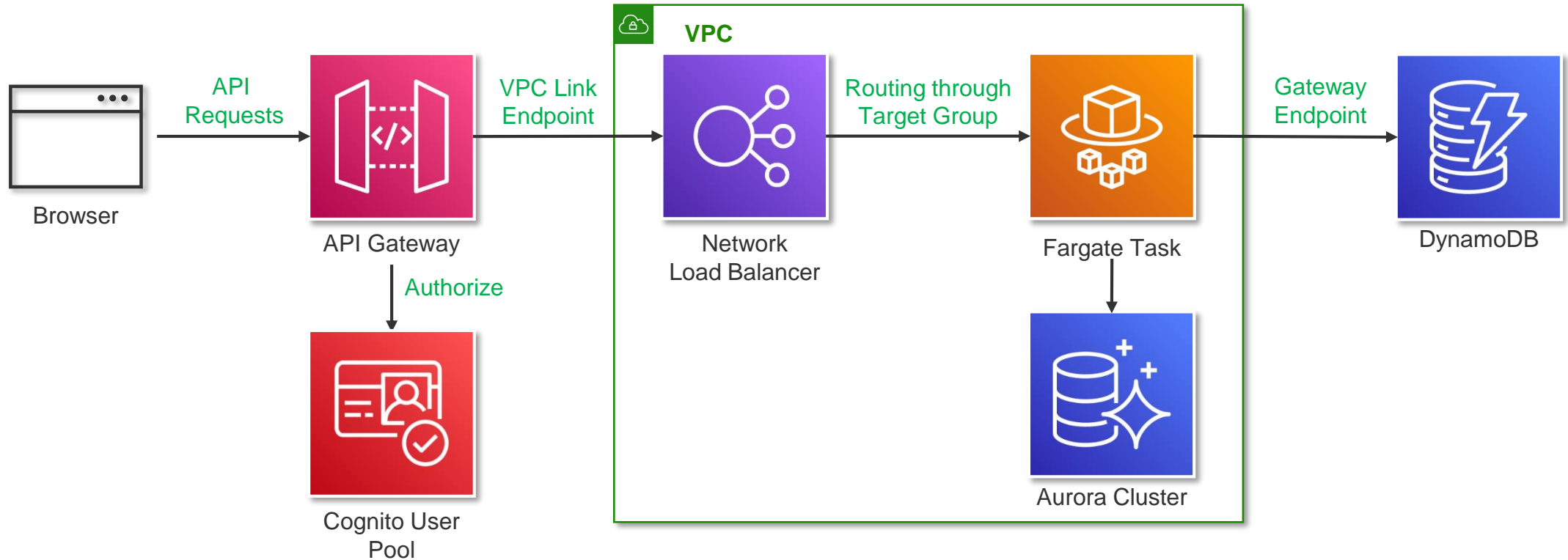


Background Processing Message Handler



Media File Processing Trigger

Containerized API Service



Pattern Samples



Prerequisites

[Install or update the latest version of the AWS CLI - AWS Command Line Interface \(amazon.com\)](#)

[Installing the AWS SAM CLI - AWS Serverless Application Model \(amazon.com\)](#)

[Getting started with the AWS CDK - AWS Cloud Development Kit \(AWS CDK\) v2 \(amazon.com\)](#)

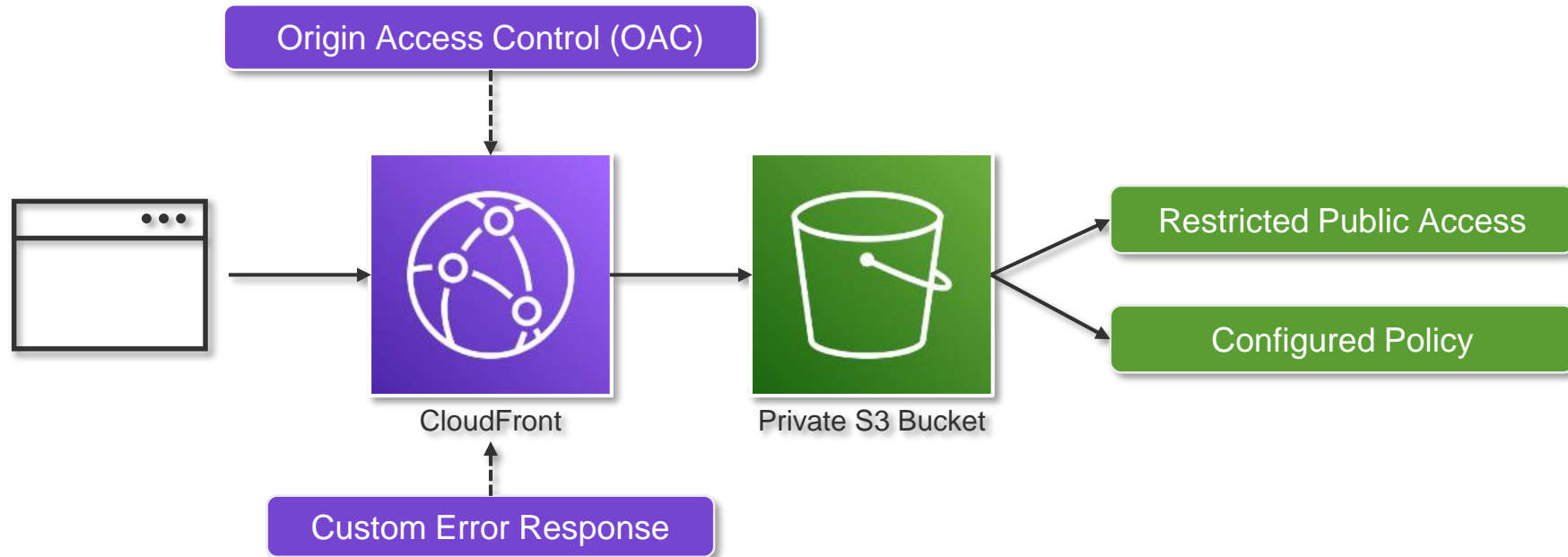


```
graph TD; A[Code Sample] --> B[luan tien/aws-serverless-samples (github.com)]
```

Code Sample

[luan tien/aws-serverless-samples \(github.com\)](#)

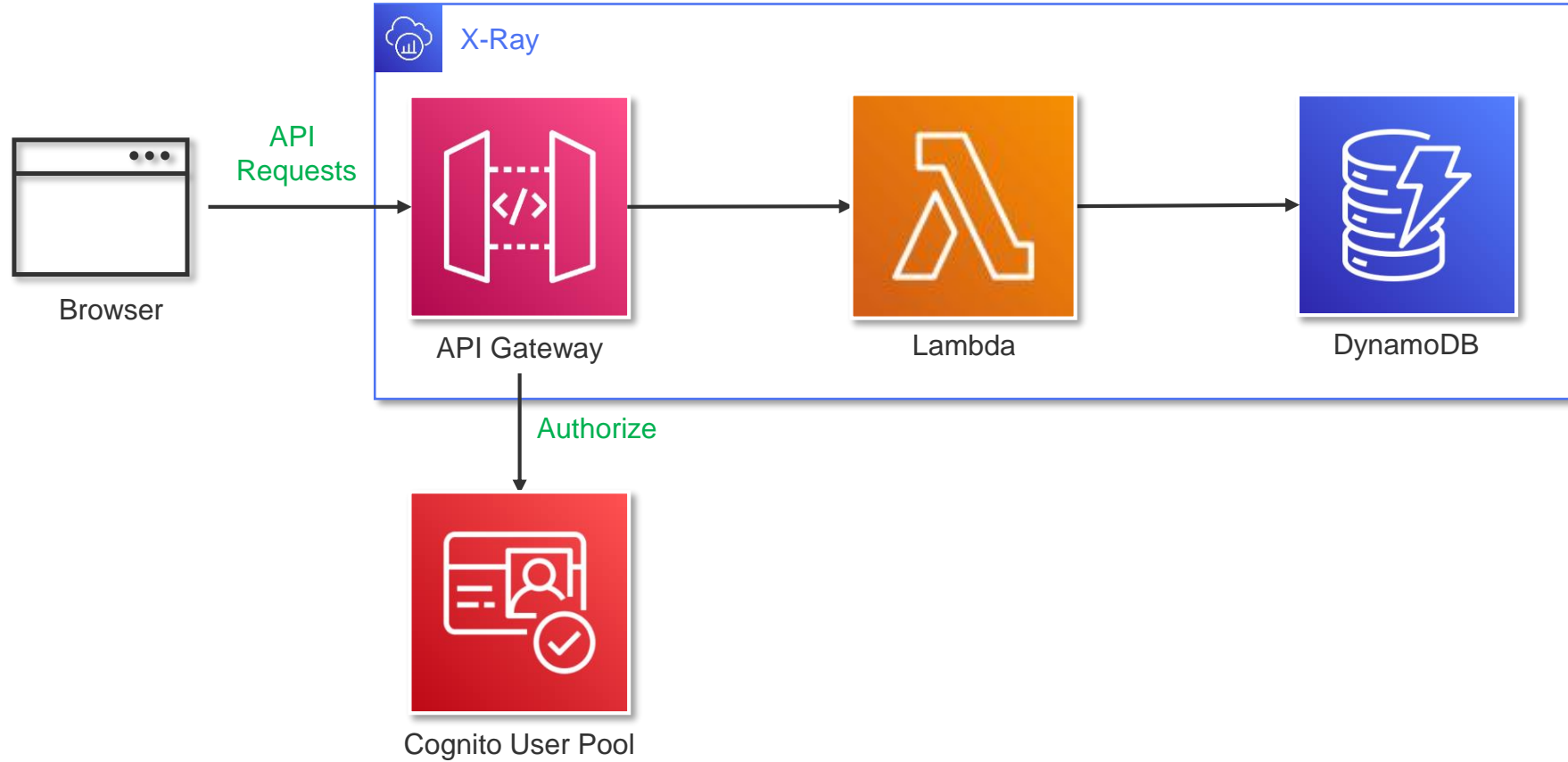
UI Service



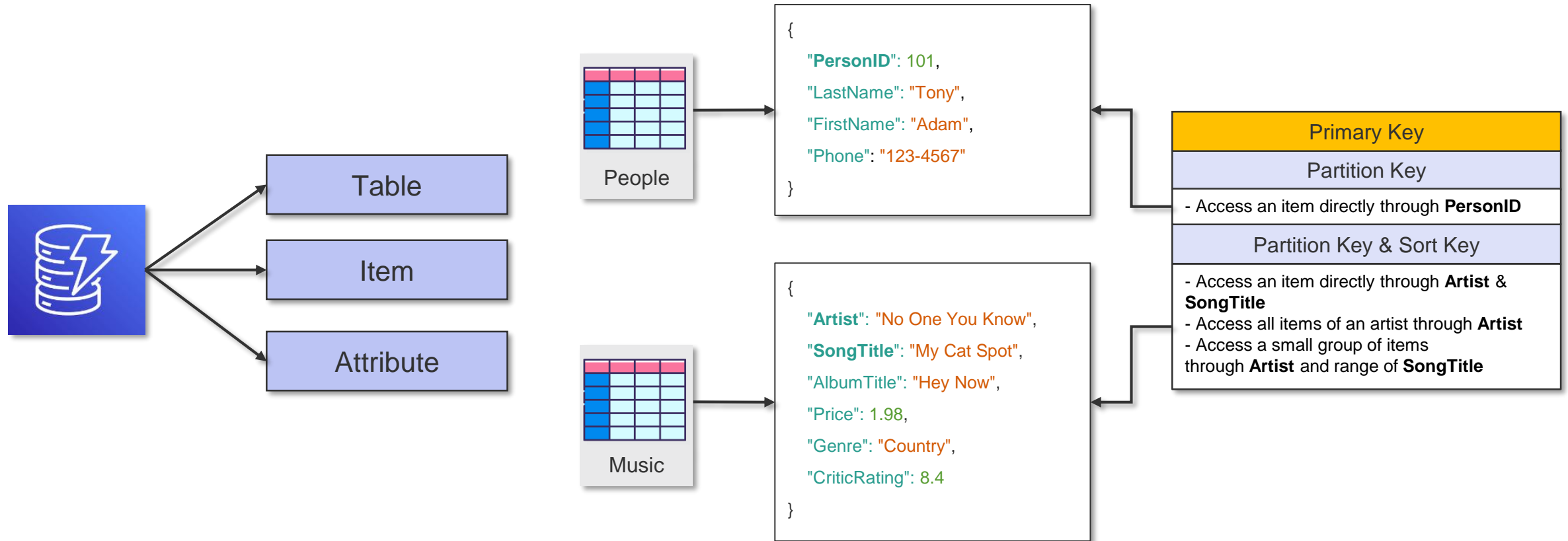
Sample Time



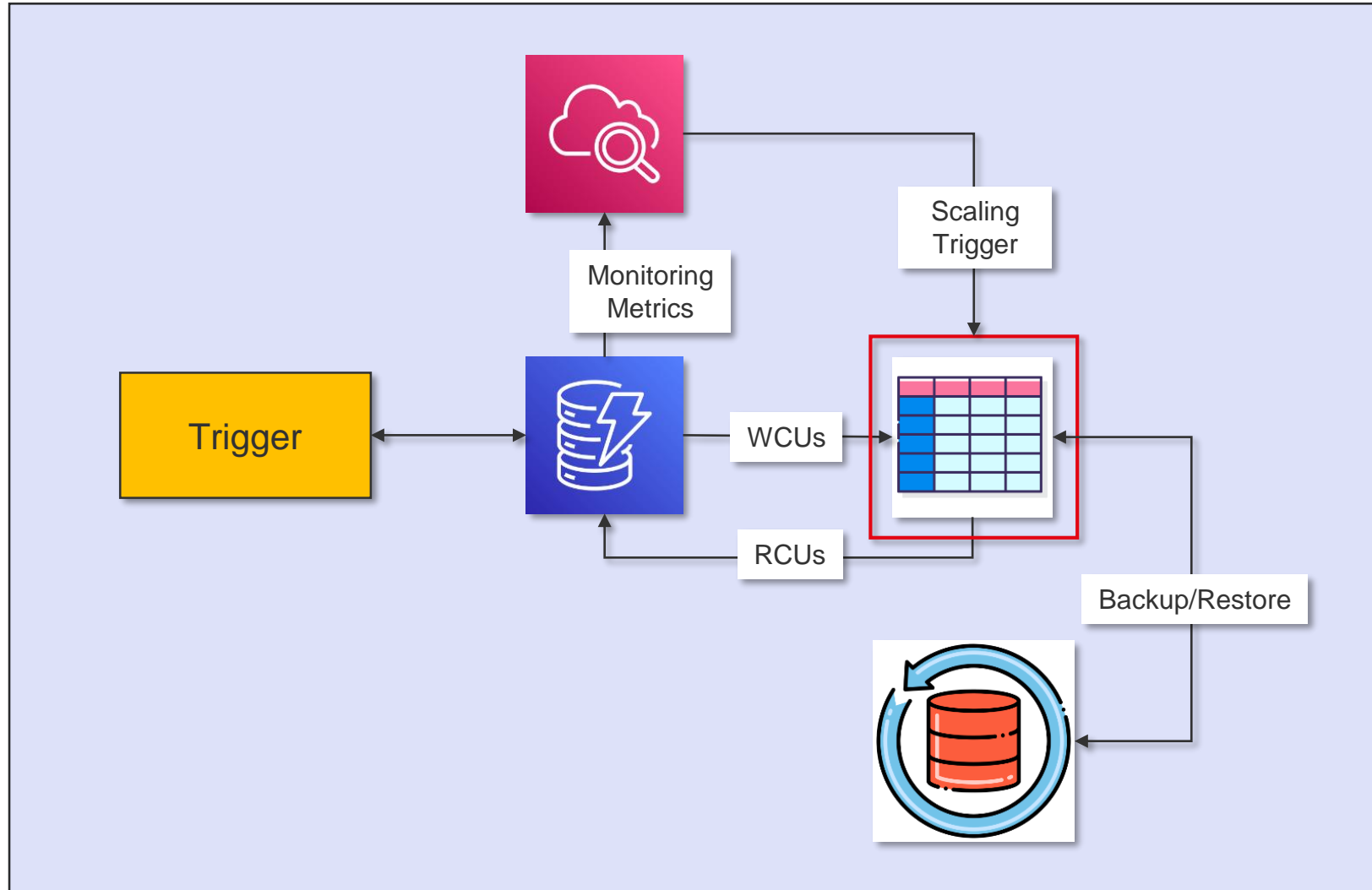
API Service with Lambda



DynamoDB



DynamoDB Operation

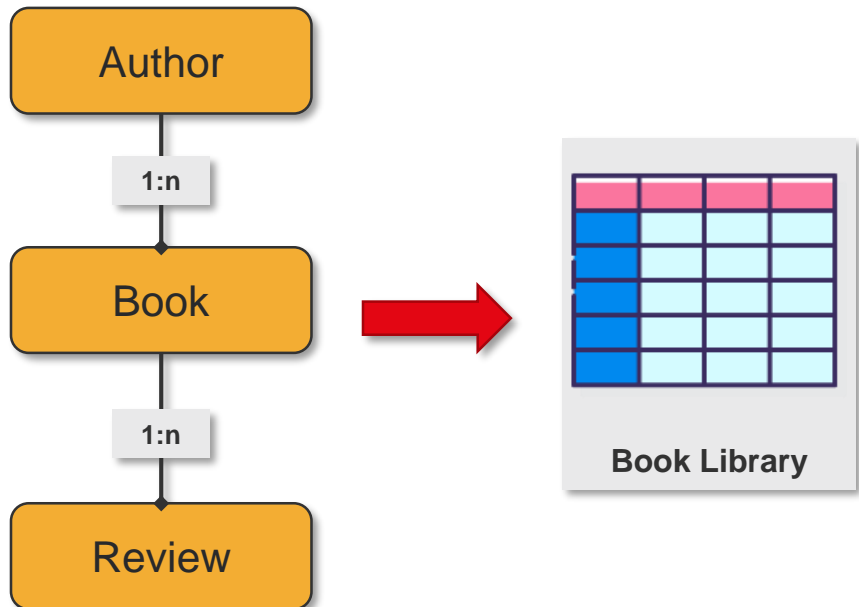


DynamoDB Pricing

Pricing	
On-demand	Provisioned
<ul style="list-style-type: none"> - Create new tables with unknown workloads. - Have unpredictable application traffic - Prefer the ease of paying for only what you use. - Risk: no threshold can lead to unexpected high payment amount. 	<ul style="list-style-type: none"> - Have predictable application traffic. - Run applications whose traffic is consistent or ramps gradually. - Can forecast capacity requirements to control costs. - Risk: capacity limit should be monitoring properly otherwise it can impact application availability and user experience.

Billing		
Features	Billing Unit	Description
Write Request	Write request unit (WCU)	<p>Write data to your table:</p> <ul style="list-style-type: none"> - Standard write request unit can write an item up to 1 KB. - Transactional write requires double more WCU than standard write. - Ex: Standard write for 3KB item is 3 WCUs. Transactional write for 3KB item is 6 WCUs.
Read Request	Read request unit (RCU)	<p>Read data from your table:</p> <ul style="list-style-type: none"> - Strongly consistent read of item up to 4 KB requires 1 RCU. - Eventually consistent read of item up to 4 KB requires 0.5 RCUs. - Transactional read of item up to 4 KB requires 2 RCUs. <p>Ex: Eventually consistent read of 8 KB item is 1 RCU. Transactional read of 8 KB item is 4 RCUs.</p>
Data Storage	GB storage per month (GB-month)	<ul style="list-style-type: none"> - Raw byte size of your data + per-item overhead of 100 bytes for indexing in table. - Ex: your table occupies 25 GB at the beginning of the month, grows to 29 GB by the end of month (30 days). Average storage in 30 days is 27 GB. We will be charged 27-25=2GB.

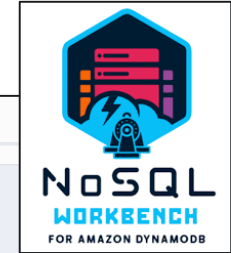
DynamoDB Table Design



BookLibrary			
Access Patterns	Table/GSI	Key Conditions	Example
Get book for a given bookId	Table	PK=bookId and SK=bookId	PK="b#12345" and SK="b#12345"
Get author for a given authorId	Table	PK=authorId and SK=authorId	PK="a#12345" and SK="a#12345"
Get all reviews for a given bookId	Table	PK=bookId and SK begins_with "r#"	PK="b#12345" and SK begins_with "r#"
Get all books for a given authorId	GSI1	PK=authorId and SK begins_with "b#"	PK="a#12345" and SK begins_with "b#"

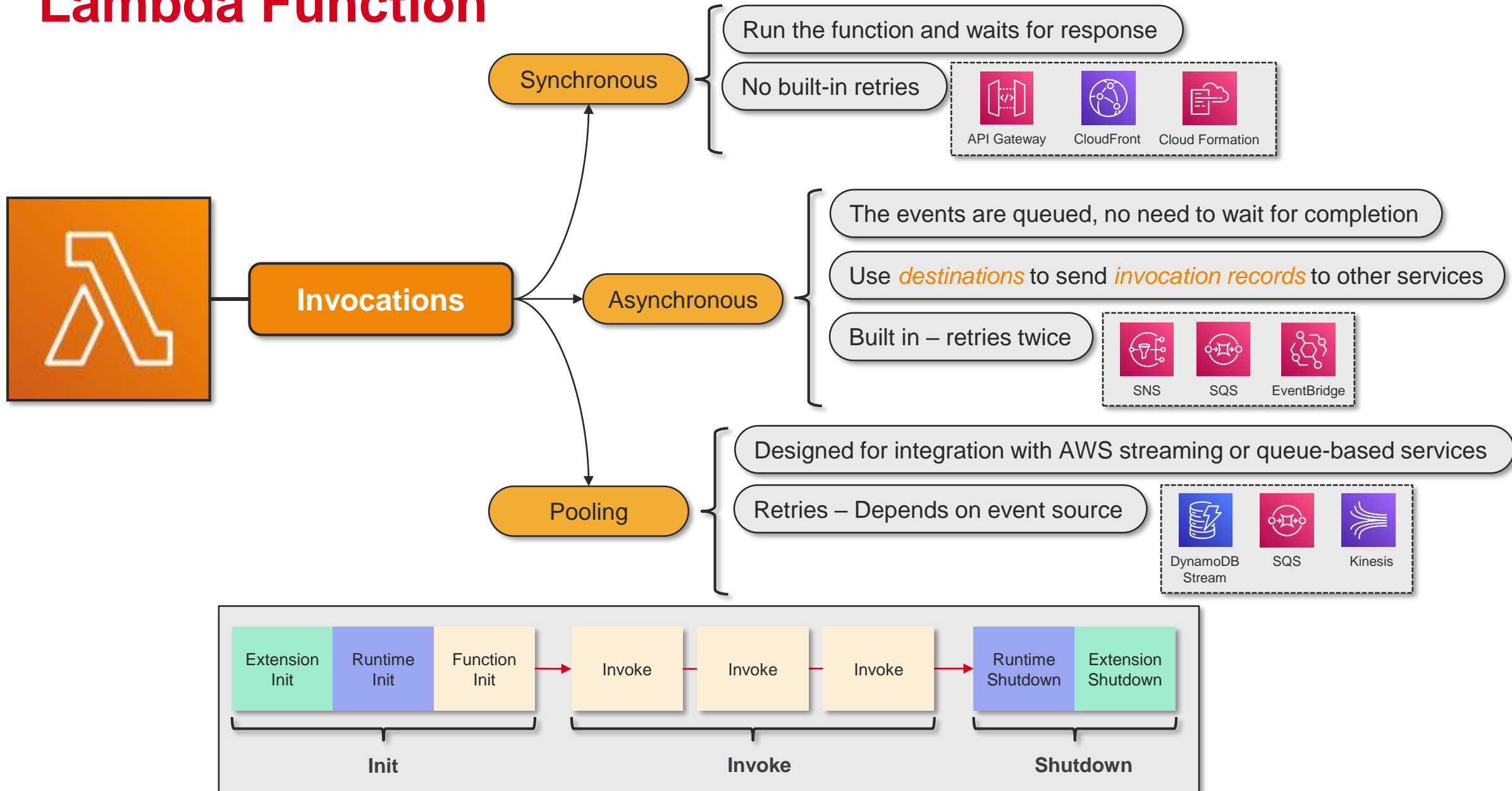
Primary Key		Attributes								
PK	SK									
b#12345	b#12345	EntityType	AuthorId (GS1-PK)	PublishedDate	Title	Description	Price	Positive	Negative	Neutral
		book	a#12345	7/10/2022	book title 1	description 1	100	0	0	0
	r#12345	EntityType	Reviewer	Title	Description					
		review	Tony Adam	Bravo book	This is a very nice book to read					
a#12345	a#12345	EntityType	Name							
		author	author 1							

DynamoDB Table Design

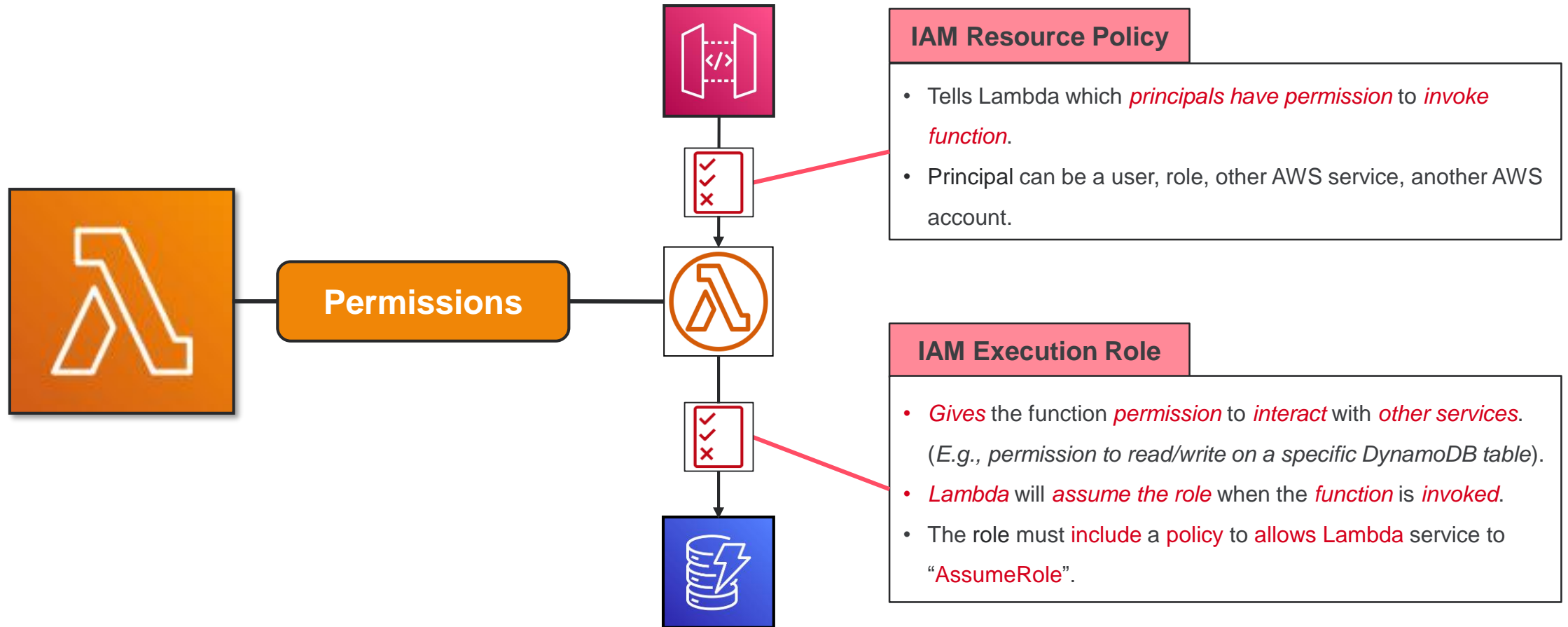


Primary key		Attributes				
Partition key: PK	Sort key: SK					
a#1	a#1	EntityType	Name			
		author	Author 1			
b#2	b#2	AuthorId	EntityType	Description	PublishedDate	Title
		a#1	book	Description 2	2019-08-01	Book 2
b#3	b#3	AuthorId	EntityType	Description	PublishedDate	Title
		a#2	book	Description 3	2018-05-22	Book 3
b#1	b#1	AuthorId	EntityType	Description	PublishedDate	Title
		a#1	book	Description 1	2018-01-01	Book 1
b#6	b#6	AuthorId	EntityType	Description	PublishedDate	Title
		a#3	book	Description 6	2020-09-03	Book 6
a#3	a#3	EntityType	Name			
		author	Author 3			
b#5	b#5	AuthorId	EntityType	Description	PublishedDate	Title
		a#3	book	Description 5	2018-11-04	Book 5
a#2	a#2	EntityType	Name			
		author	Author 2			
b#4	b#4	AuthorId	EntityType	Description	PublishedDate	Title
		a#2	book	Description 4	2021-07-20	Book 4

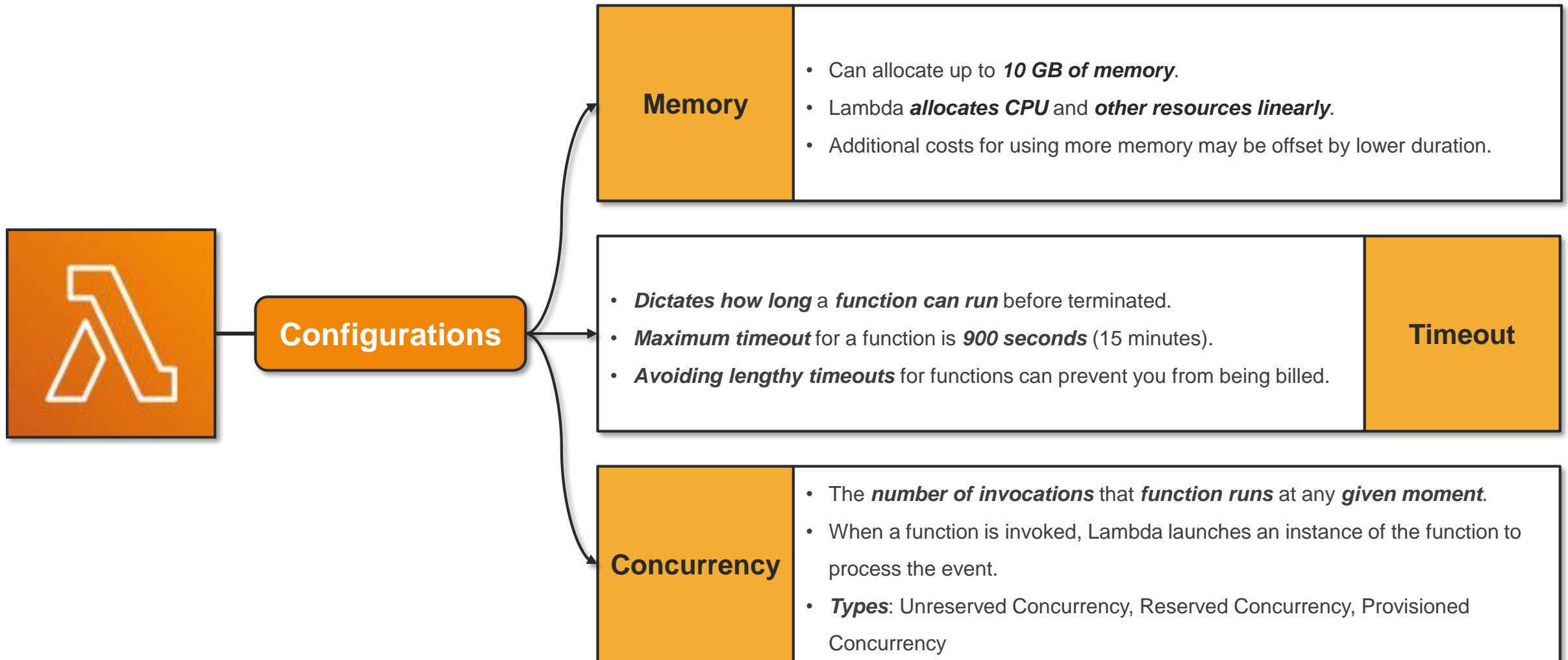
Lambda Function



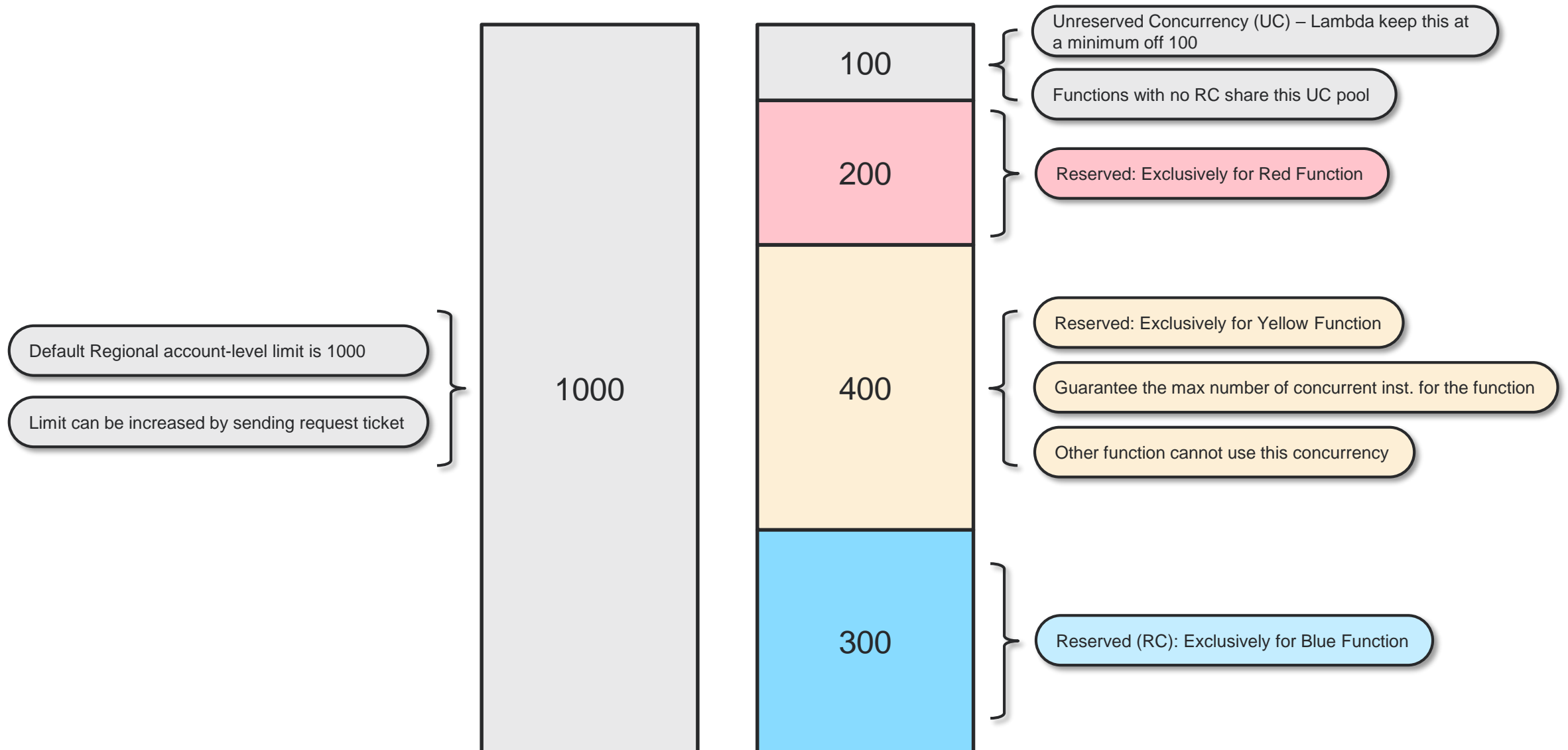
Lambda Function Permission



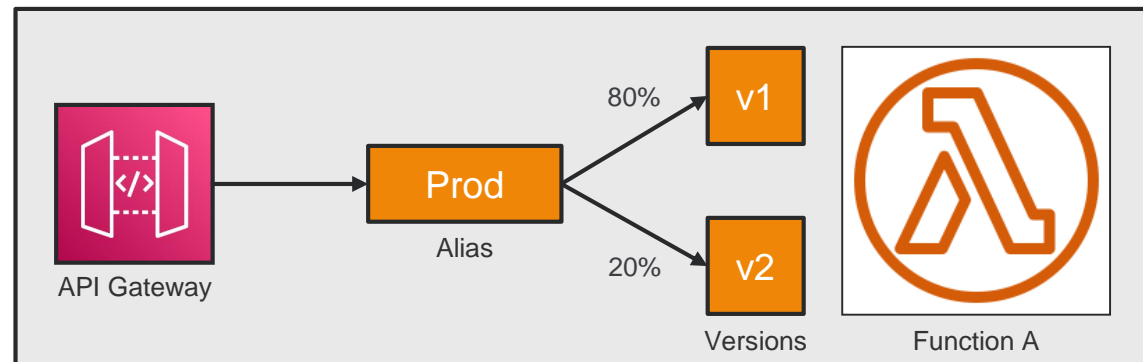
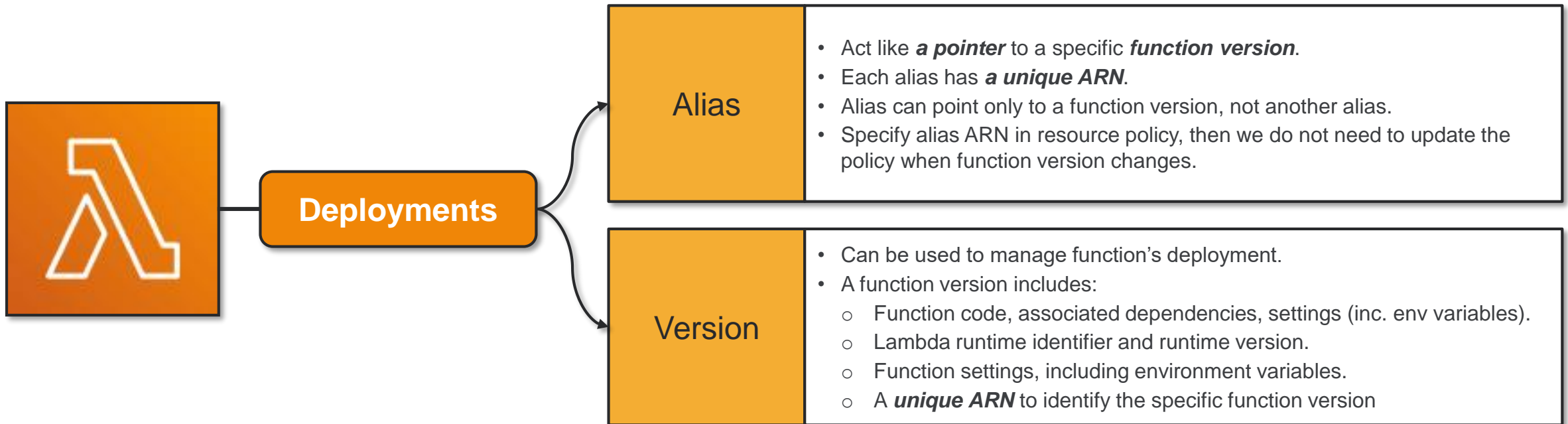
Lambda Configuration



Lambda Concurrency Samples



Lambda Deployment



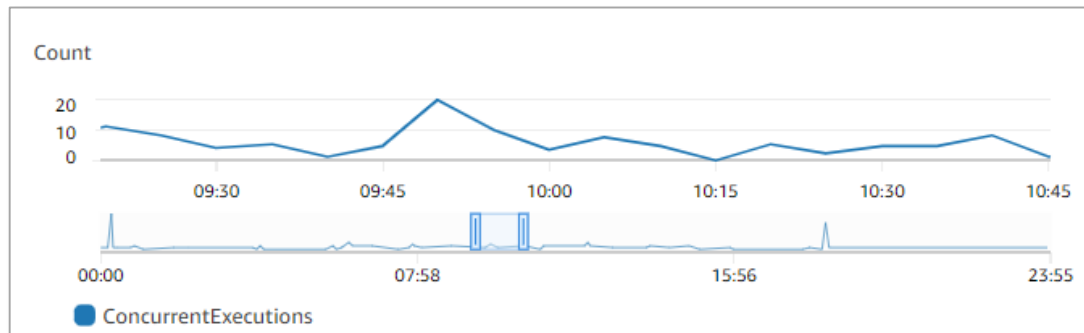
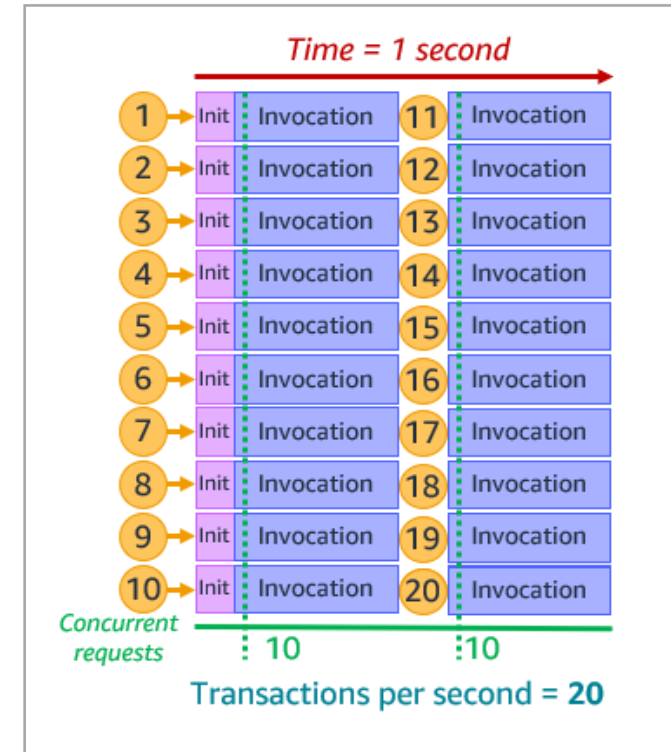
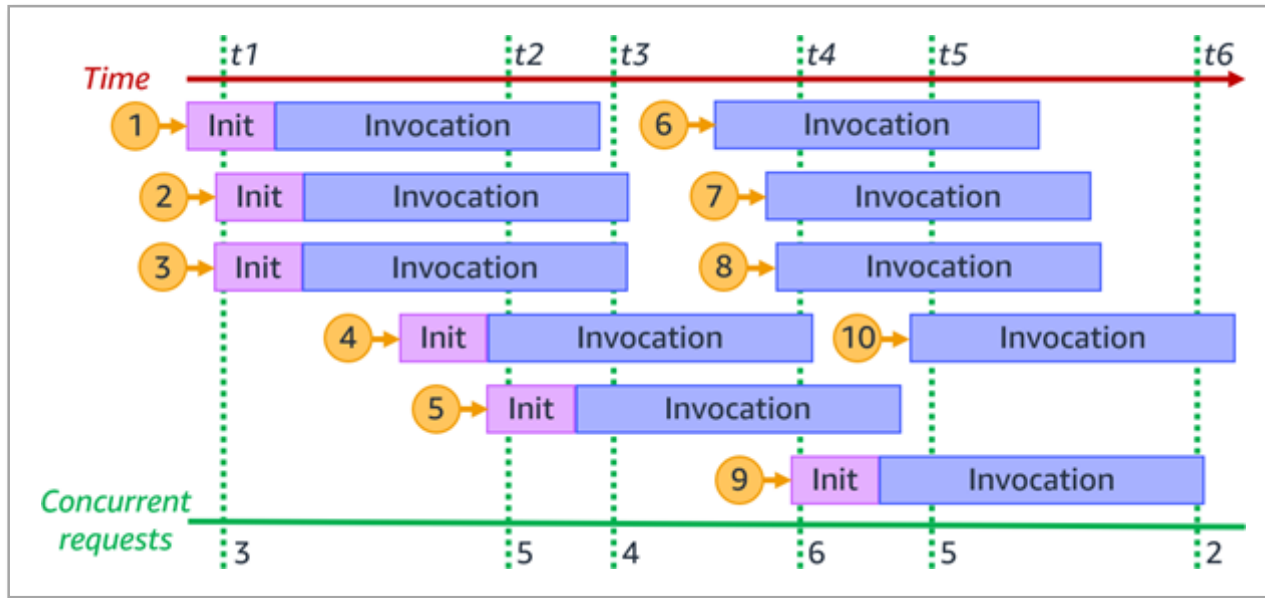
Lambda Scaling and Quota

Service Quotas			
Service Quotas > AWS services > AWS Lambda			
AWS Lambda			
Service quotas Request quota increase			
Filter by...			
Service quota	Applied quota value	AWS default quota value	Adjustable
<input type="radio"/> Concurrent executions	Not available	1,000 / Second	Yes
<input type="radio"/> Elastic network interfaces per VPC	Not available	250	Yes
<input type="radio"/> Function and layer storage	Not available	75 Gigabytes	Yes
<input checked="" type="radio"/> Asynchronous payload	Not available	256 Kilobytes	No
<input checked="" type="radio"/> Burst concurrency	Not available	3,000	No
<input checked="" type="radio"/> Deployment package size (direct upload)	Not available	50 Megabytes	No
<input checked="" type="radio"/> Deployment package size (unzipped)	Not available	250 Megabytes	No
<input checked="" type="radio"/> Environment variable size	Not available	4 Kilobytes	No
<input checked="" type="radio"/> File descriptors	Not available	1,024	No
<input checked="" type="radio"/> Function memory maximum	Not available	3,008 Megabytes	No
<input checked="" type="radio"/> Function memory minimum	Not available	128 Megabytes	No
<input checked="" type="radio"/> Function timeout	Not available	900	No
<input checked="" type="radio"/> Processes and threads	Not available	1,024	No
<input checked="" type="radio"/> Synchronous payload	Not available	6 Megabytes	No
<input checked="" type="radio"/> Temporary storage	Not available	512 Megabytes	No

Up to **10,000** per **region**

- **3000** – US West (Oregon), US East (N. Virginia), Europe (Ireland)
- **1000** – Asia Pacific (Tokyo), Europe (Frankfurt), US East (Ohio)
- **500** – Other Regions

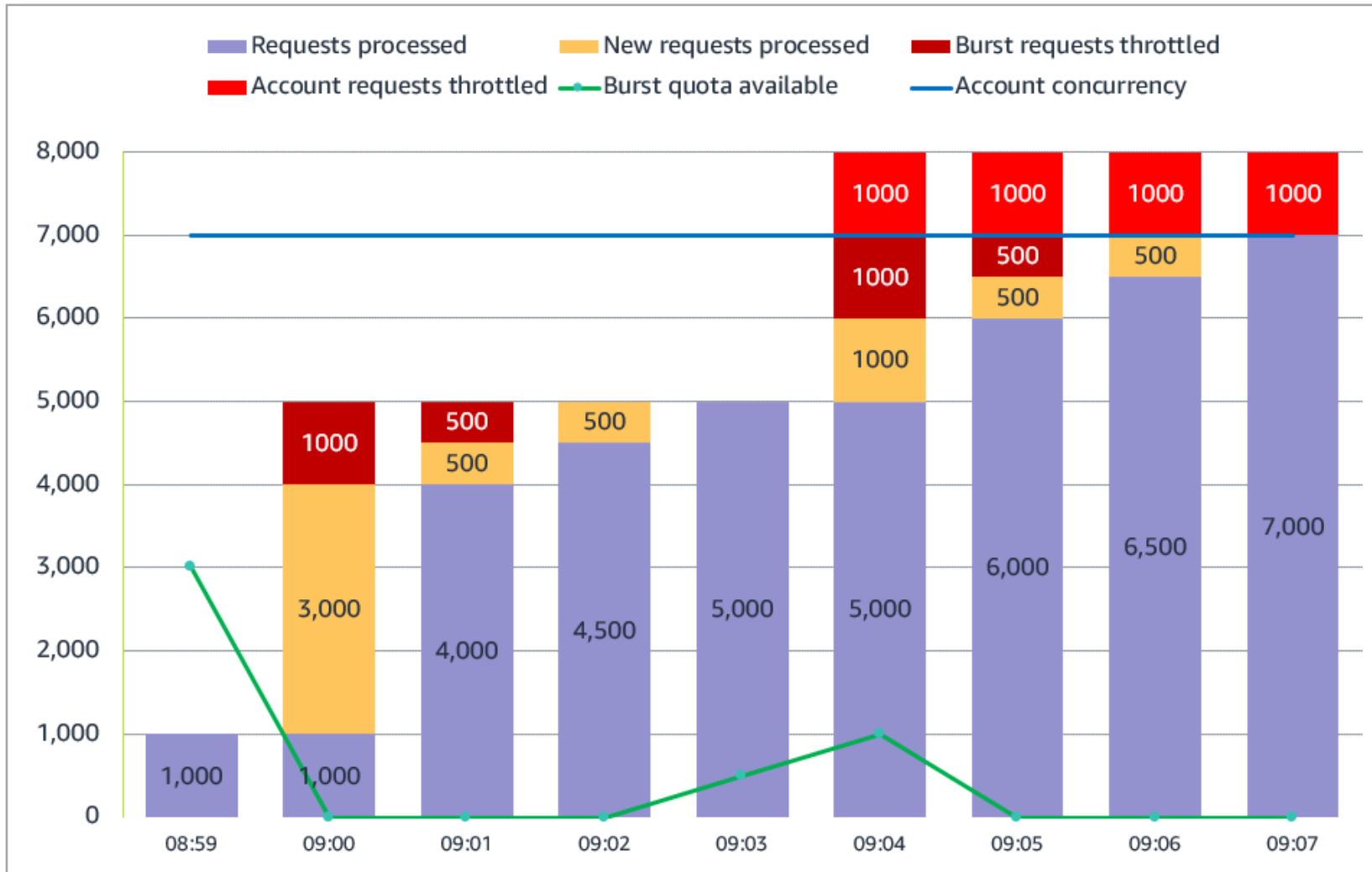
Lambda Scaling and Quota



$RequestsPerSecond \times AvgDurationInSeconds = concurrent\ requests$

$100\ requests/second \times 0.5\ sec = 50\ concurrent\ requests$

Lambda Scaling and Quota



[Understanding AWS Lambda scaling and throughput | AWS Compute Blog \(amazon.com\)](https://aws.amazon.com/blogs/compute/understanding-aws-lambda-scaling-and-throughput/)

08:59

- Using **1,000** concurrent **runtime environments**.
- Lambda invocation takes **250ms**.
- **Transactions per second (TPS)** are **4,000**.

09:00

- Large traffic with **5,000 sustained requests**.
- **1000 requests executed** by existing runtime environments.
- Lambda uses the **3,000 available burst concurrency** to create new environments to handle the additional load.
- **1,000 requests are throttled** as there is not enough burst concurrency to handle all 5,000 requests.
- **TPS are 16,000**.

09:01

- Lambda scales by another **500 concurrent invocations per minute**.
- **500 requests** are still throttled.
- App can now handle **4,500 concurrent requests**.

09:02

- Lambda scales by another **500 concurrent invocations per minute**.
- No requests are throttled.
- App can now handle all **5,000 concurrent requests**.

09:03

- The application continues to handle the **sustained 5000 requests**.
- The **burst concurrency quota** rises to **500**.

09:04

- The application continues to handle the **sustained 5000 requests**.
- The **burst concurrency quota** rises to **500**.

09:05

- Lambda scales by another **500 concurrent requests**.
- The application can now handle **6,500 requests**.
- **500 requests are throttled** as there is **not enough burst concurrency**.
- **1,000 requests are still throttled** as the **account concurrency quota** has been **reached**.

Lambda Pricing

Allocated memory: 512 MB

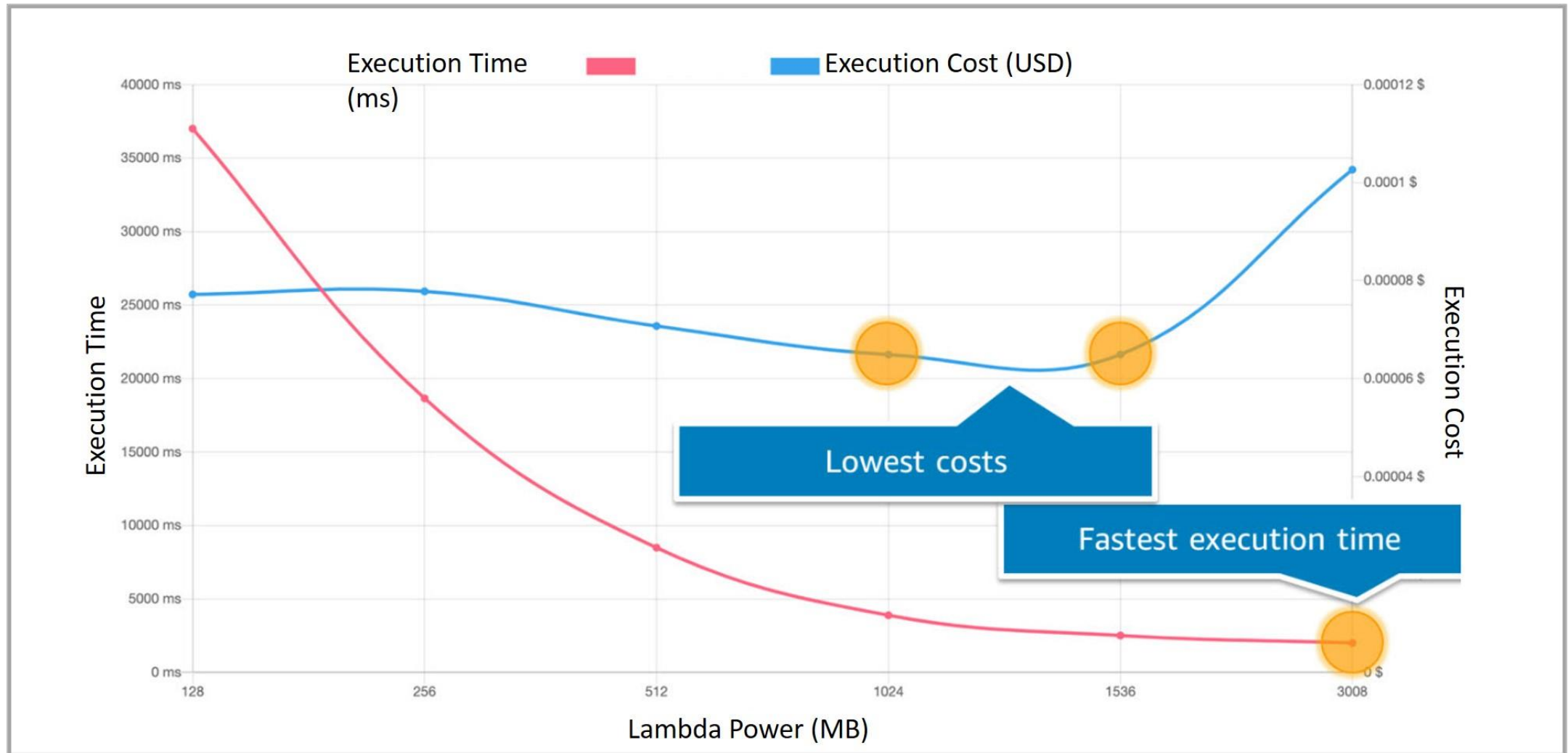
No. Invocations: 20,000 times/month

Execution duration: 1sec

AWS LAMBDA	
GB-sec 20,000 * 512/1024	10,000 GB-sec
Compute Charged 10,000 * 0.00001667	\$0.1667
Requests Charged (20,000/1 Mil) * 0.2 per Mil	\$0.004
TOTAL	\$0.1707

EC2 (on-demand t2.nano)	
\$0.0081*24*30	\$5.832
TOTAL	\$5.832

Lambda Pricing

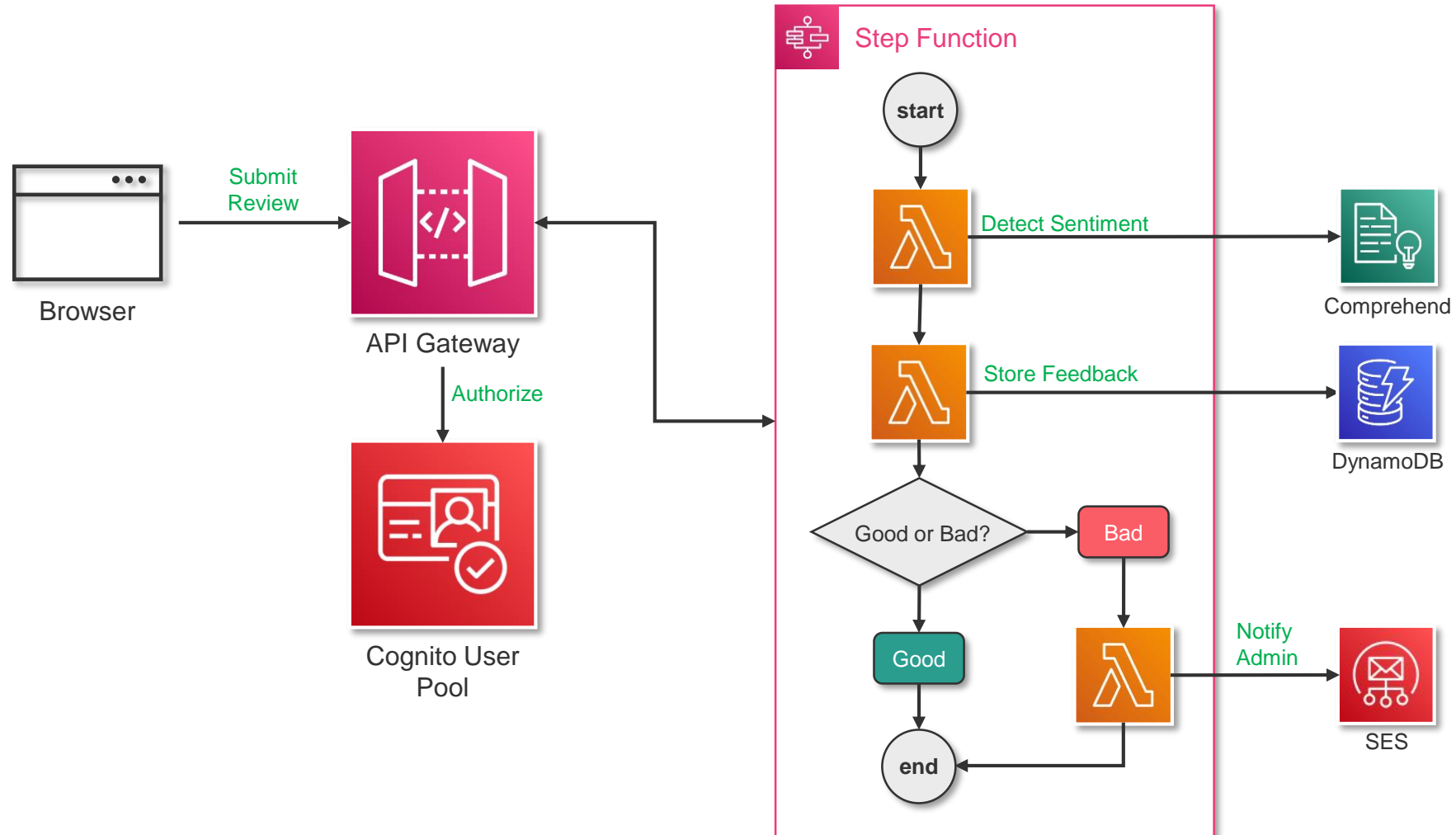


Lambda Power Tuning provides a comparison of cost vs. speed of your Lambda function


Sample Time



Workflow Service (Sync)



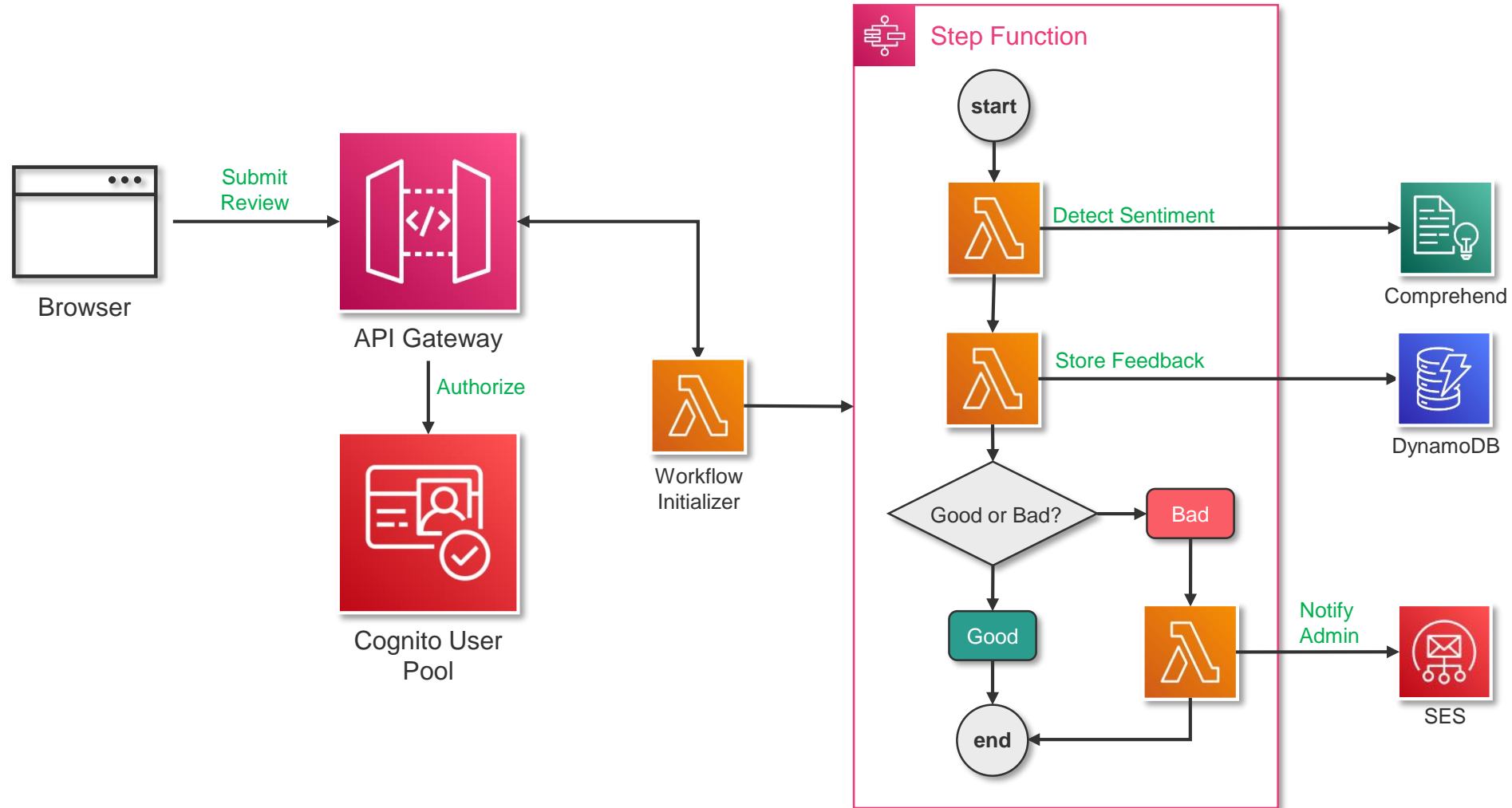
State Machine Type

Standard vs Express Workflows		
	Standard Workflows	Express Workflows: Synchronous and Asynchronous
Maximum duration	1 year.	5 minutes.
Supported execution start rate	Over 2,000 per second	Over 100,000 per second
Supported state transition rate	Over 4,000 per second per account	Nearly unlimited
Pricing 	Priced per state transition. A state transition is counted each time a step in your execution is completed.	Priced by the number of executions you run, their duration, and memory consumption.
Execution history	Executions can be listed and described with Step Functions APIs, and visually debugged through the console. They can also be inspected in CloudWatch Logs by enabling logging on your state machine.	Unlimited execution history, that is, as many execution history entries are maintained as you can generate within a 5-minute period. Further, executions can be inspected in CloudWatch Logs by enabling logging on your state machine.
Execution semantics	Exactly-once workflow execution.	<i>Asynchronous Express Workflows:</i> At-least-once workflow execution. <i>Synchronous Express Workflows:</i> At-most-once workflow execution.
Service integrations	Supports all service integrations and patterns.	Supports all service integrations. Note Express Workflows do not support Job-run (.sync) or Callback (.waitForTaskToken) service integration patterns.
Step Functions activities	Supports Step Functions activities.	Does not support Step Functions activities.

State Machine Type

Standard Workflows	Asynchronous Express Workflows	Synchronous Express Workflows
Exactly-once workflow execution	At-least-once workflow execution	At-most-once workflow execution
<p>The execution state is internally persisted on every state transition. To guarantee that only one workflow with the same name can run, Step Functions will return an idempotent response when you start a Standard Workflow with the same name as an already running workflow. In this case, Step Functions will not start a new workflow. When the workflow completes, Step Functions will respond with an exception. After 90 days, the workflow data will be removed, and the name can then be reused.</p>	<p>No internally persisted state for workflow progress. If you attempt to start an Express Workflow with the same name more than once, each attempt causes a workflow to start concurrently. In rare cases, the internal state of a workflow can be lost, and the workflow will be automatically restarted from beginning. You should ensure your state machine logic is idempotent and should not be affected adversely by multiple concurrent executions of the same input.</p>	<p>After a workflow starts, Step Functions will wait and returns the result as part of the API response. If service exceptions occur, Step Functions will not restart from the beginning. You should ensure your state machine logic is idempotent and should not be affected adversely by multiple concurrent executions of the same input.</p>

Workflow Service (Async)



Sample Time



Thank you