# AWS re:Invent

**NOV. 28 – DEC. 2, 2022 | LAS VEGAS, NV**

SVS305

# Understanding AWS Lambda concurrency

Brian Zambrano (he/him)

Sr. Specialist Solutions Architect, Serverless
Amazon Web Services

Justin Pirtle (he/him)

Principal Solutions Architect
Amazon Web Services

**"We want a serverless architecture that supports 10,000 requests per second."**

- The business

"While keeping costs under control,

with good performance,

and integration with our relational database."

# AWS Lambda concurrency affects all of this
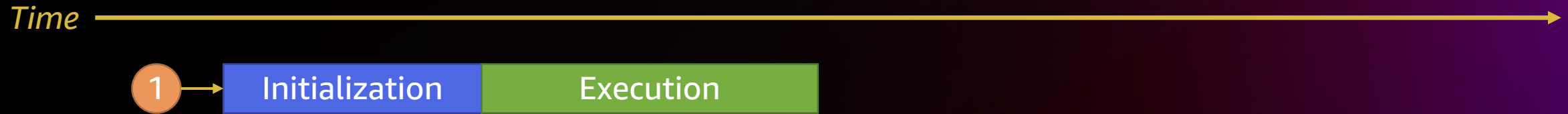
# Agenda

What we will cover

- AWS Lambda concurrency and concurrency controls

- Applying concurrency to a sample Serverless API architecture

- Best practices for scaling synchronous and asynchronous workloads

What we will not cover

- What is serverless?

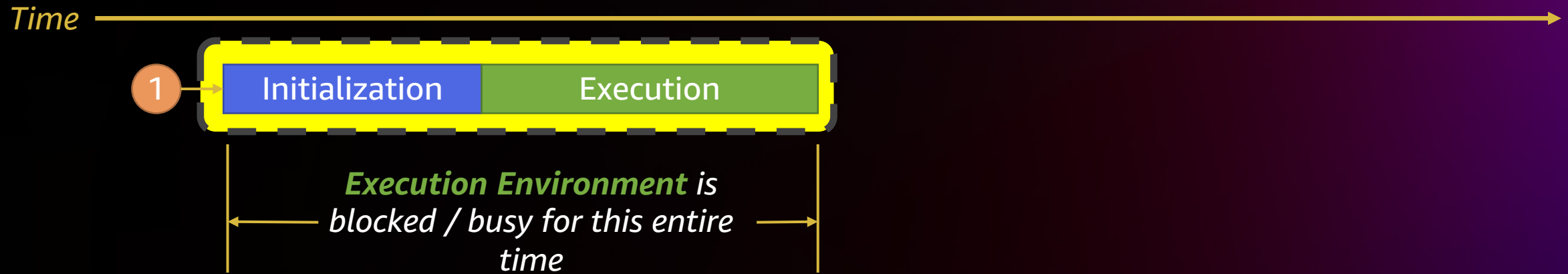- What are Amazon API Gateway, Amazon DynamoDB, AWS Lambda?
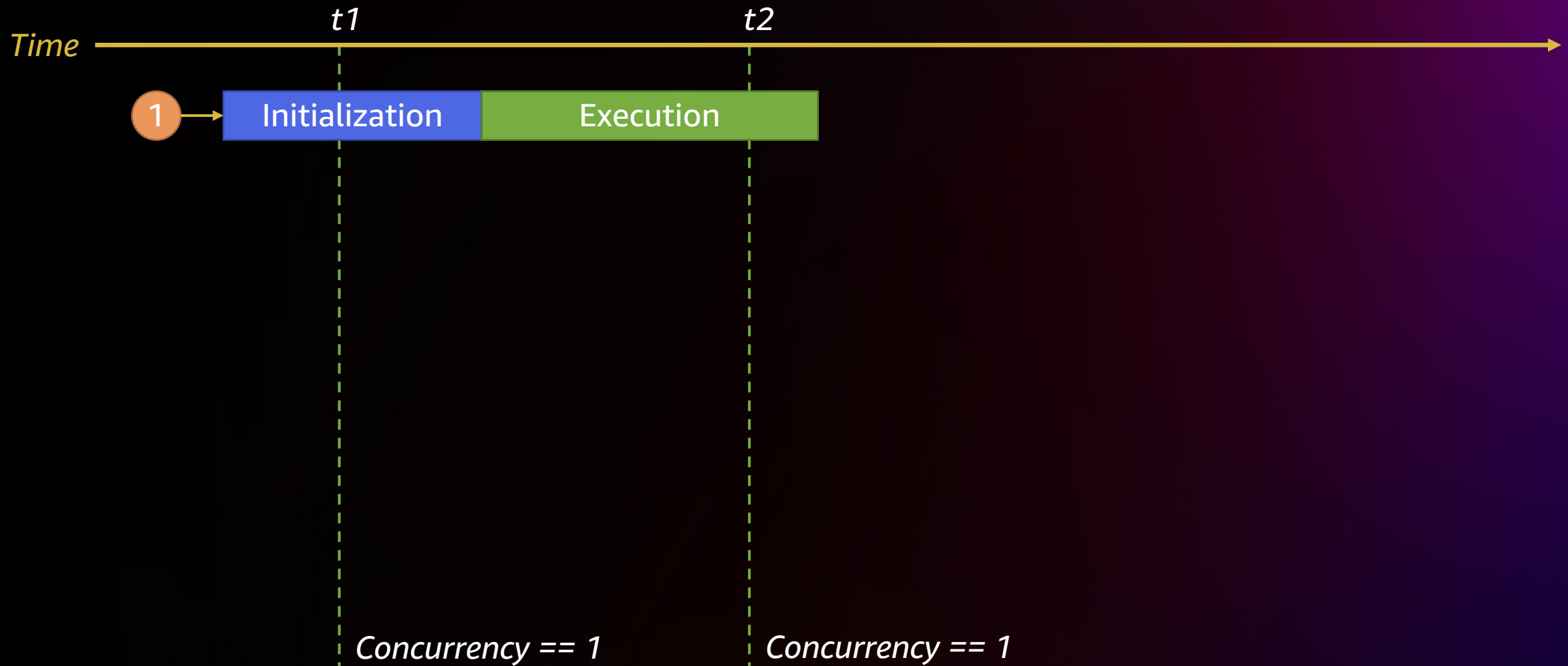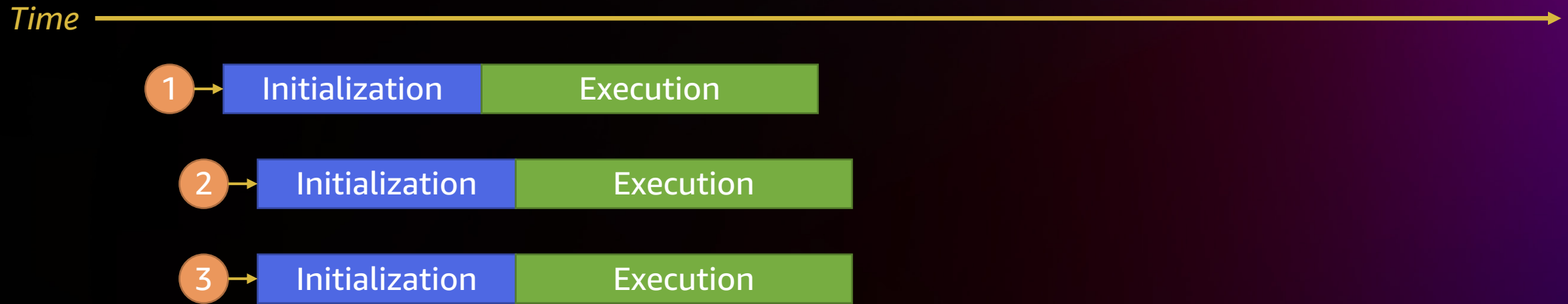
# Lambda concurrency

Time →

1 → | Initialization | Execution |

# Lambda concurrency

Time

1 | Initialization | Execution |

# Lambda concurrency

Time



**1** Initialization | Execution

*Execution Environment is blocked / busy for this entire time*

# Lambda concurrency

Time →

t1                    t2

**1** → Initialization | Execution

Concurrency == 1        Concurrency == 1

# Lambda concurrency

Time

1 → | Initialization | Execution |

2 → | Initialization | Execution |
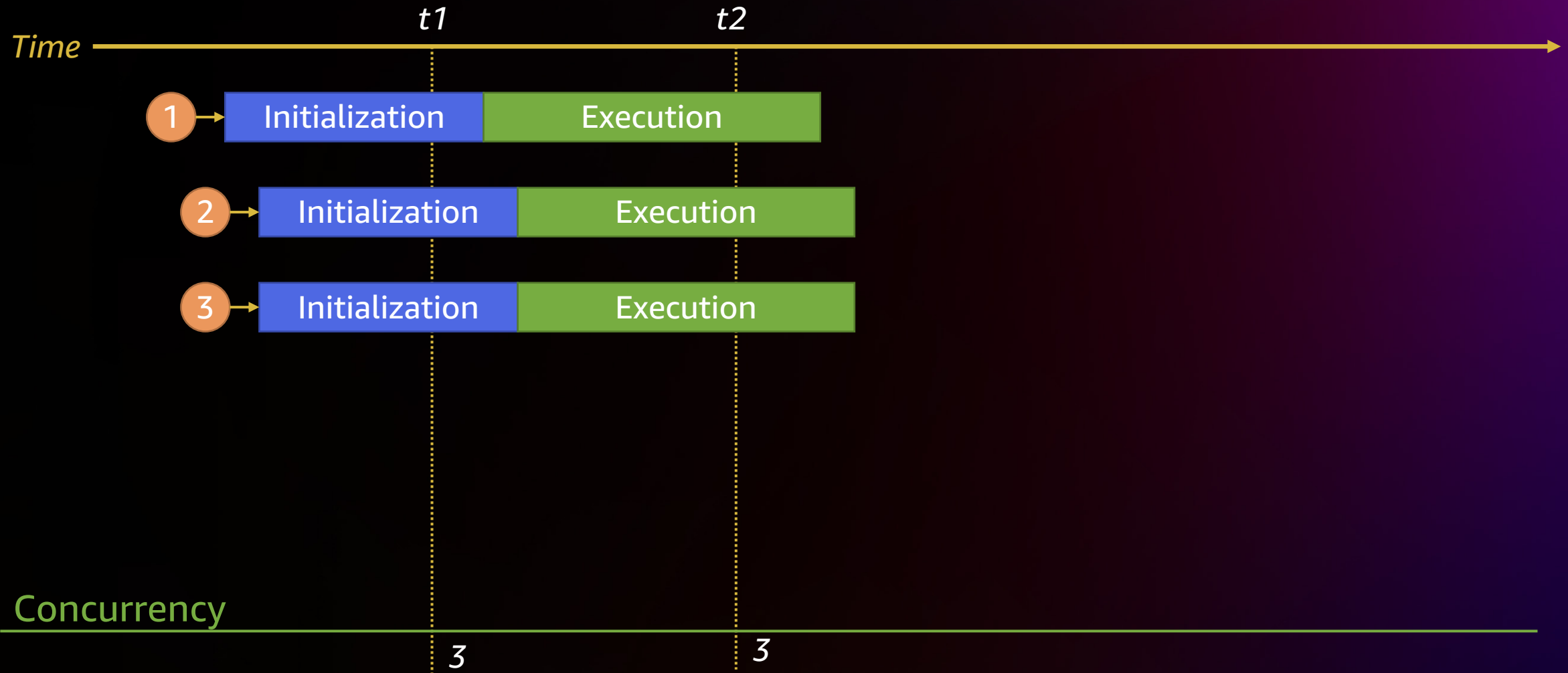
3 → | Initialization | Execution |

# Lambda concurrency

# Lambda concurrency
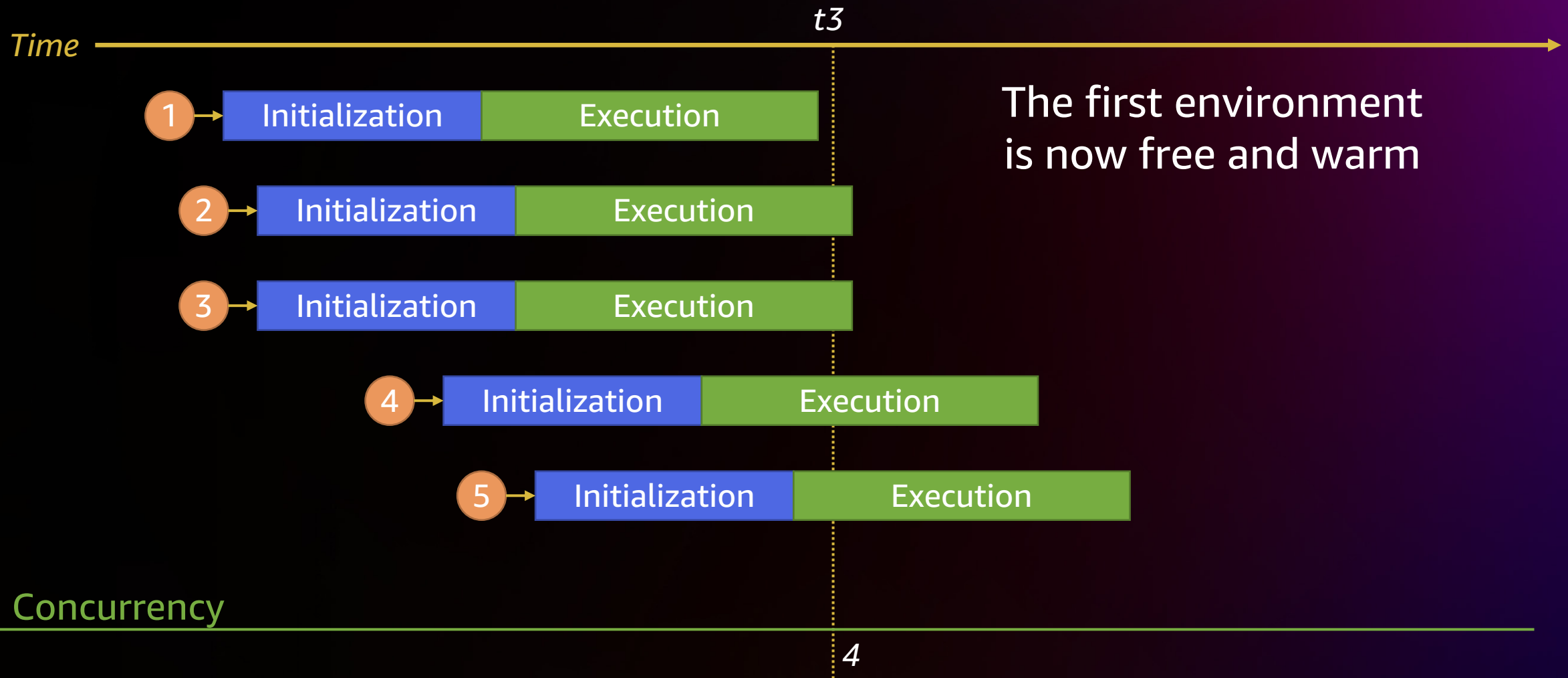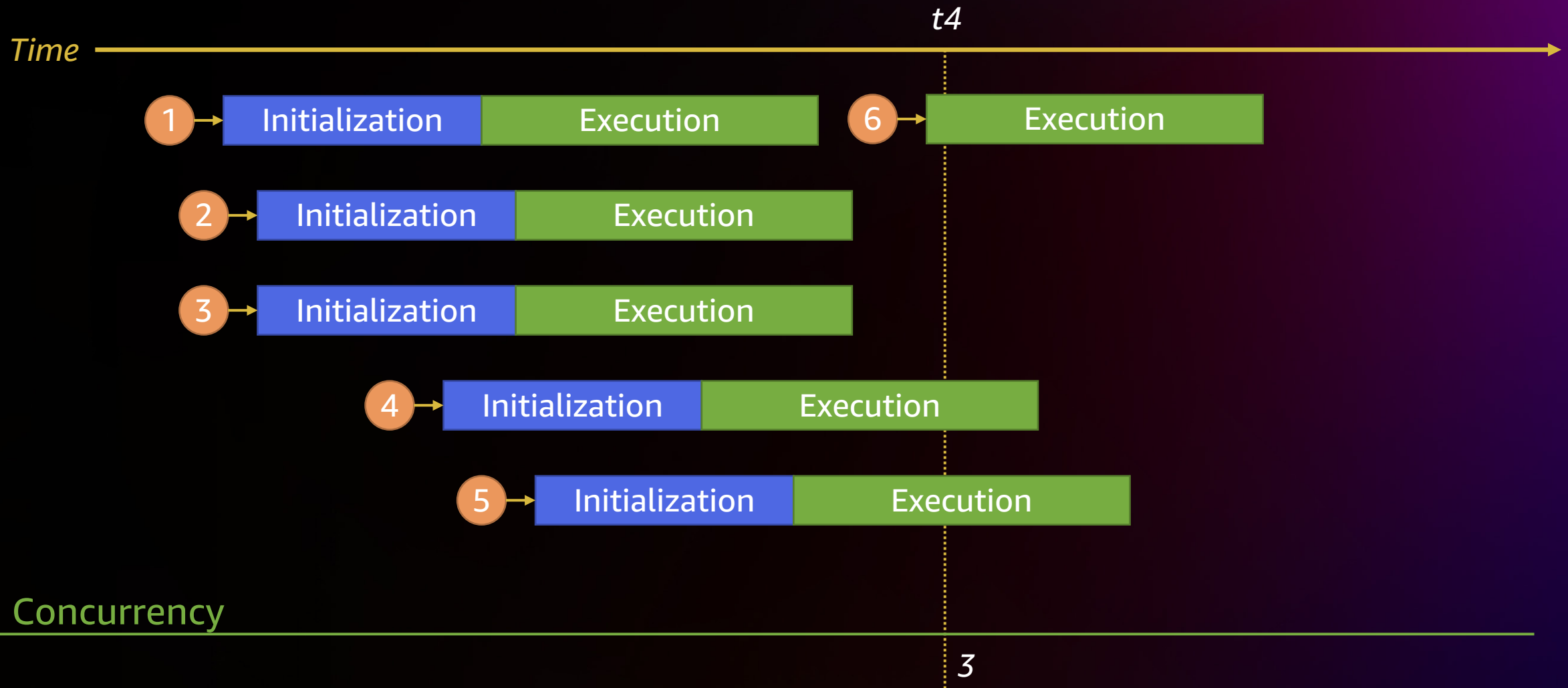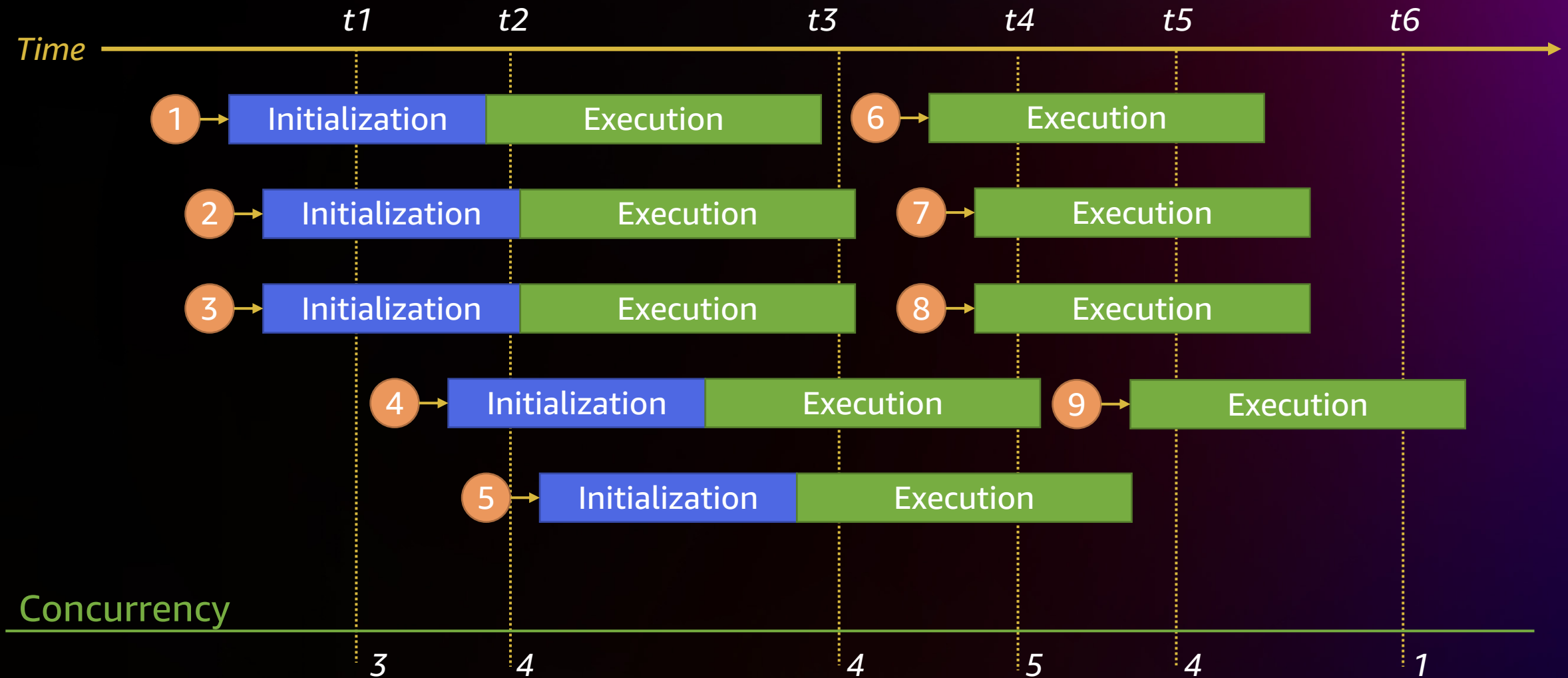
# Lambda concurrency

# Lambda concurrency

# Concurrency is a point-in-time measurement

# Concurrency is a point-in-time measurement

and not the same as requests per second.

# How to calculate concurrency requirements
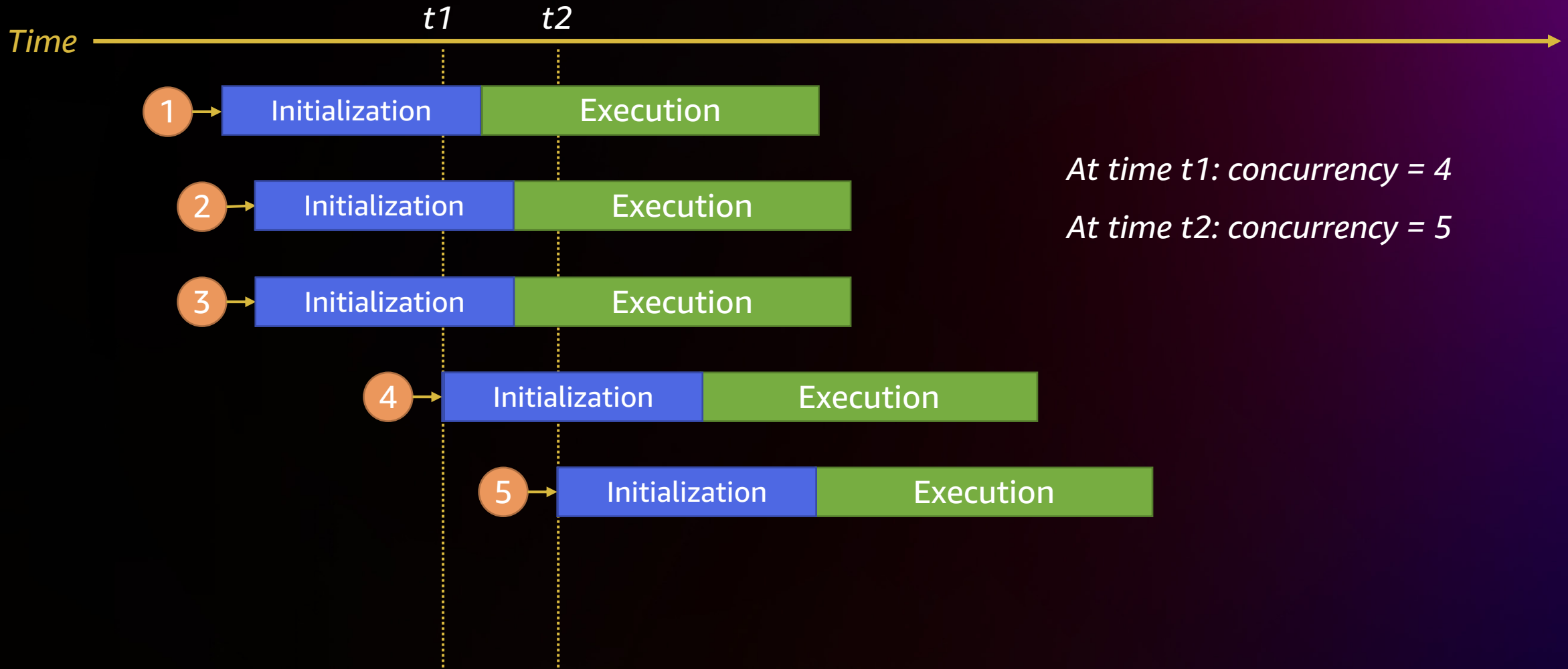
# Lambda concurrency

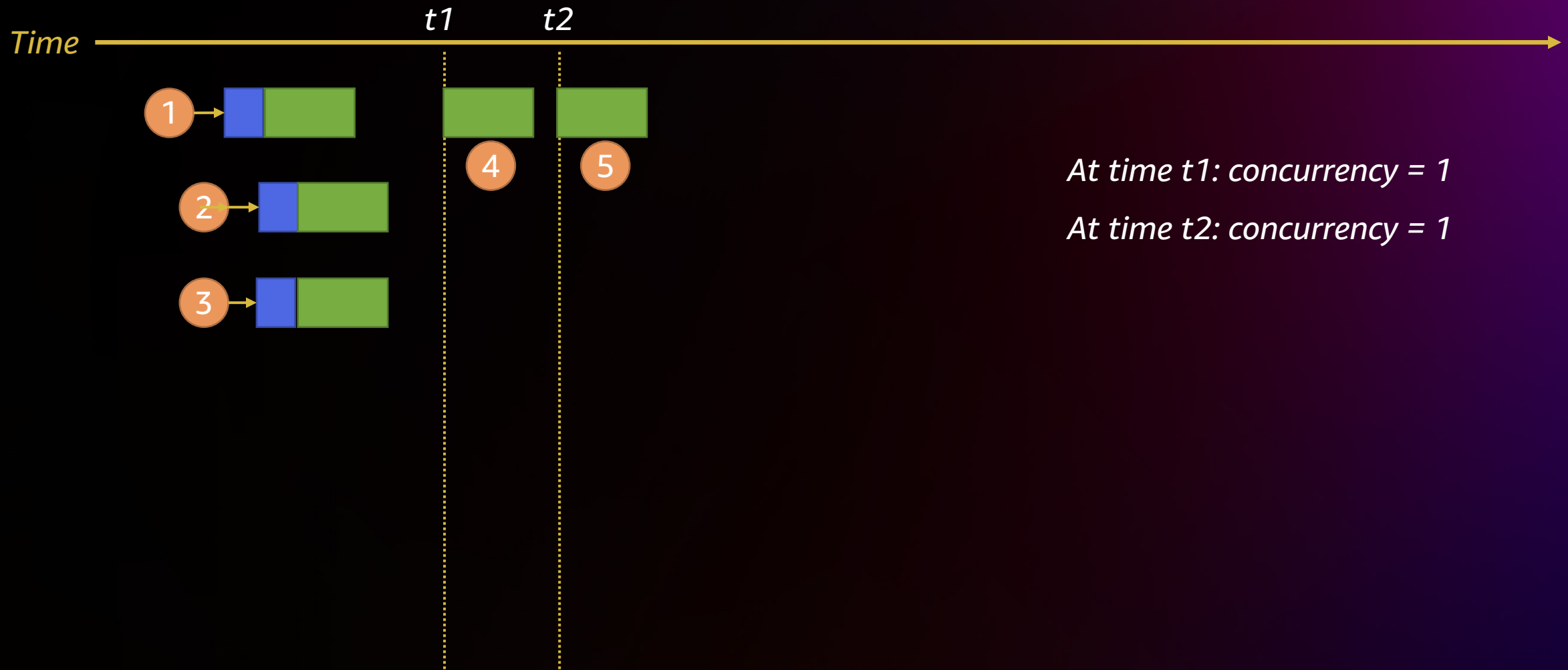**Lambda concurrency = requests per second (RPS) x duration in seconds**

Examples

- 10,000 RPS x 1,000 ms = 10,000 concurrency

- 10,000 RPS x 500 ms = 5,000 concurrency

- 10,000 RPS x 100 ms = 1,000 concurrency

**Long-running functions require more concurrency!**
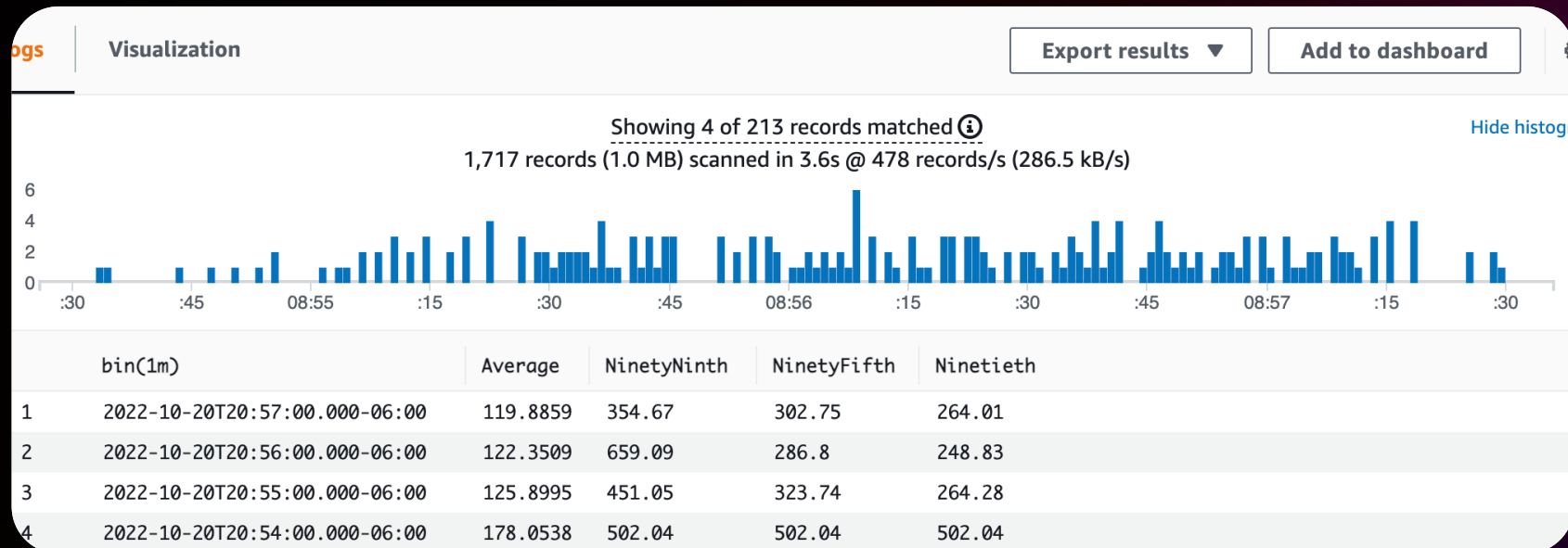
# Decreasing concurrency

# Decreasing concurrency



At time t1: concurrency = 1

At time t2: concurrency = 1

# Measuring warm starts

```
filter @type = "REPORT" and @message not like /(?i)(Init Duration)/
| stats
    avg(@duration) as Average,
    pct(@duration, 99) as NinetyNinth,
    pct(@duration, 95) as NinetyFifth,
    pct(@duration, 90) as Ninetieth
by bin(1m)
```
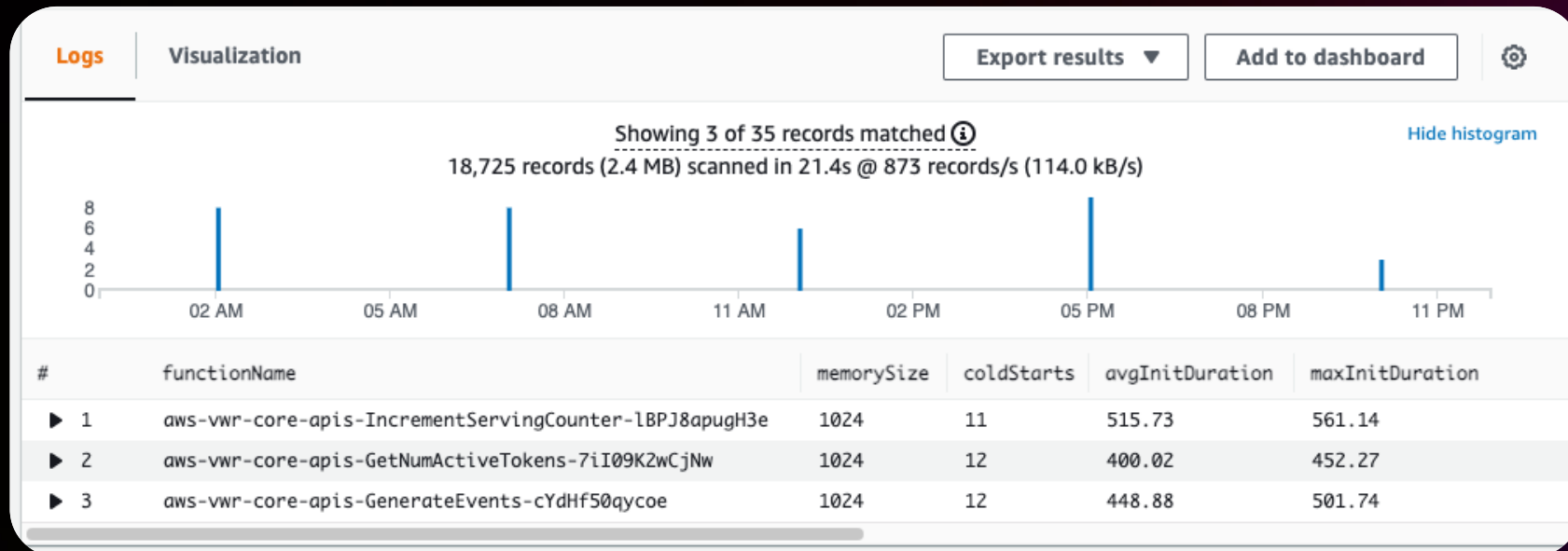


| Visualization | | Export results ▼ | Add to dashboard |
|---|---|---|---|

Showing 4 of 213 records matched ⓘ
1,717 records (1.0 MB) scanned in 3.6s @ 478 records/s (286.5 kB/s)

| | bin(1m) | Average | NinetyNinth | NinetyFifth | Ninetieth |
|---|---|---|---|---|---|
| 1 | 2022-10-20T20:57:00.000-06:00 | 119.8859 | 354.67 | 302.75 | 264.01 |
| 2 | 2022-10-20T20:56:00.000-06:00 | 122.3509 | 659.09 | 286.8 | 248.83 |
| 3 | 2022-10-20T20:55:00.000-06:00 | 125.8995 | 451.05 | 323.74 | 264.28 |
| 4 | 2022-10-20T20:54:00.000-06:00 | 178.0538 | 502.04 | 502.04 | 502.04 |

**Lambda optimizations: https://s12d.com/Lambda-Optimizations-2021**

# Measuring cold starts

```
filter @type="REPORT" and @message like /(?i)(Init Duration)/
| parse @message /^REPORT.*Init Duration: (?<initDuration>.*) ms.*/
| parse @log /^.*\/aws\/lambda\/(?<functionName>.*)/
| fields @memorySize / 1000000 as memorySize
| stats count() as coldStarts, avg(initDuration) as avgInitDuration, max(initDuration) as maxInitDuration
by functionName, memorySize
```
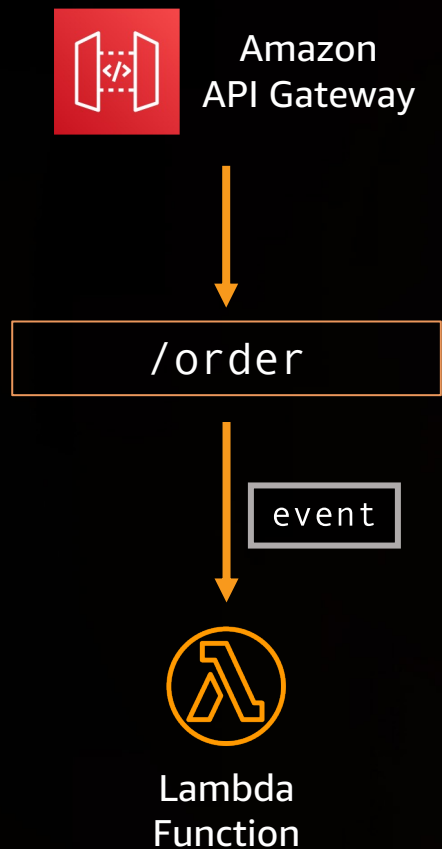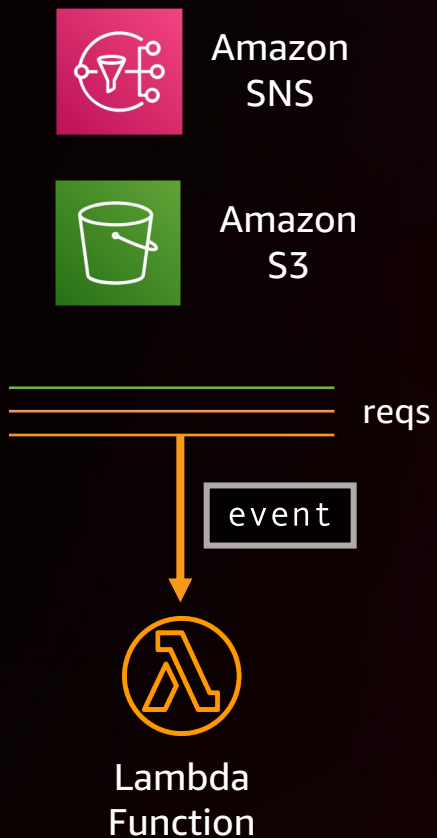


**Lambda optimizations: https://s12d.com/Lambda-Optimizations-2021**

# Applied concurrency: Lambda invocation models

# Lambda invocation methods

**Synchronous**

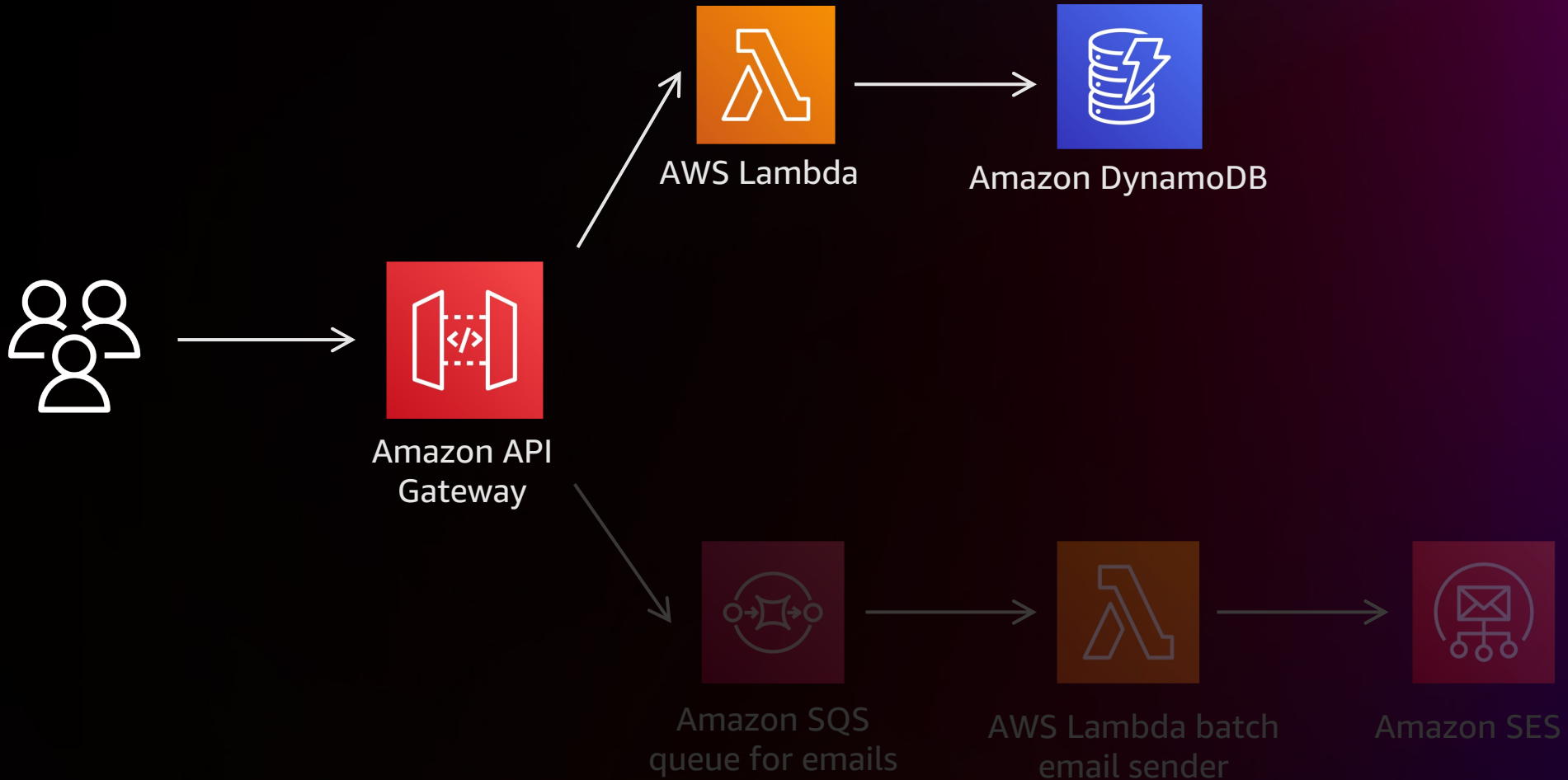Amazon
API Gateway

/order

event

Lambda
Function

**Asynchronous**

Amazon
SNS

Amazon
S3

reqs

event

Lambda
Function

**Poll-Based**

Amazon
SQS

Amazon
Kinesis

Changes

AWS Lambda
Service

event

Lambda
Function

# Today's application



AWS Lambda → Amazon DynamoDB

Amazon API Gateway

Amazon SQS queue for emails → AWS Lambda batch email sender → Amazon SES

# Synchronous invocations

# AWS Lambda function scaling quotas

## Account concurrency quota

Maximum concurrency in a  given Region
across all  functions in an account

1,000 per Region by default*

This can be increased to customer needs.

# AWS Lambda function scaling quotas

## Burst concurrency quota

Maximum increase in concurrency for an initial burst of traffic

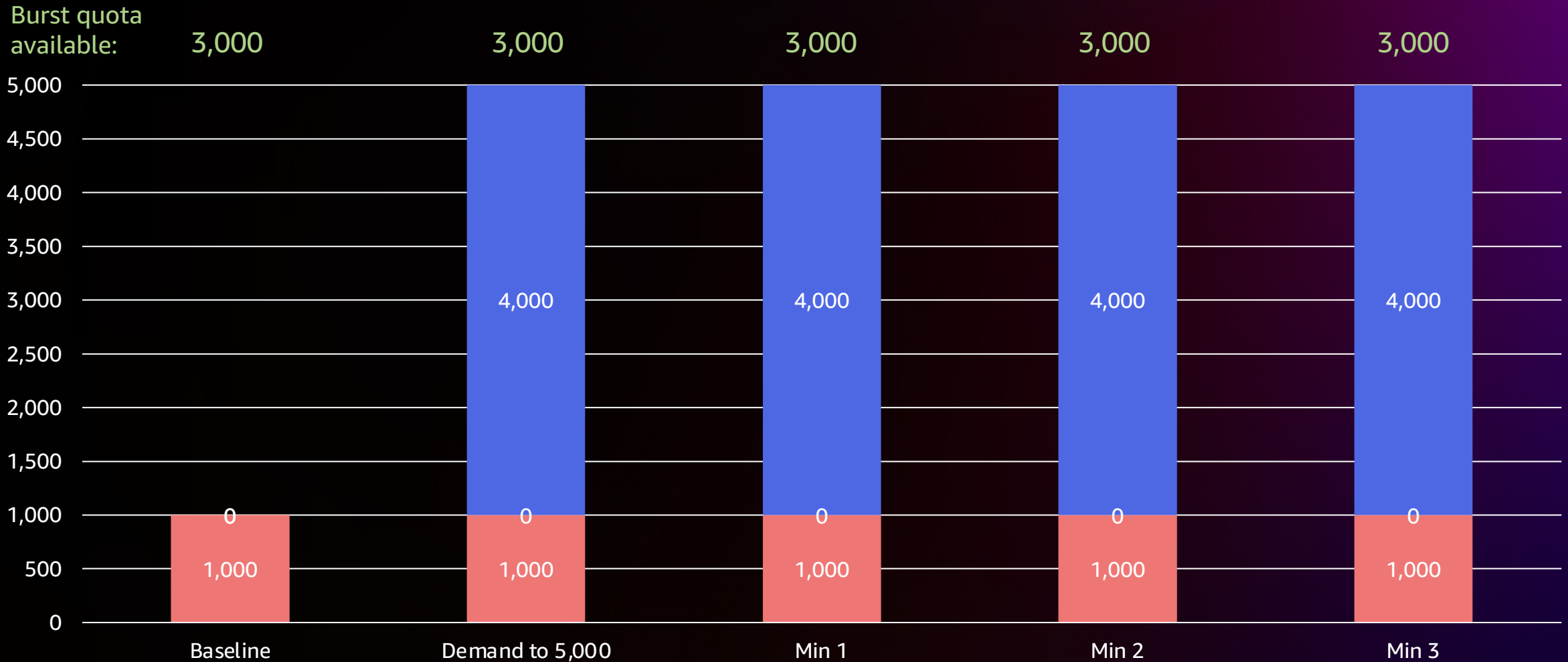**3000** in Oregon, N. Virginia and Ireland

**1000** in Tokyo, Frankfurt and Ohio

**500** in all other regions

After the initial burst, your function concurrency can scale by an additional **500 instances each minute**

# Lambda bursting and ramp up
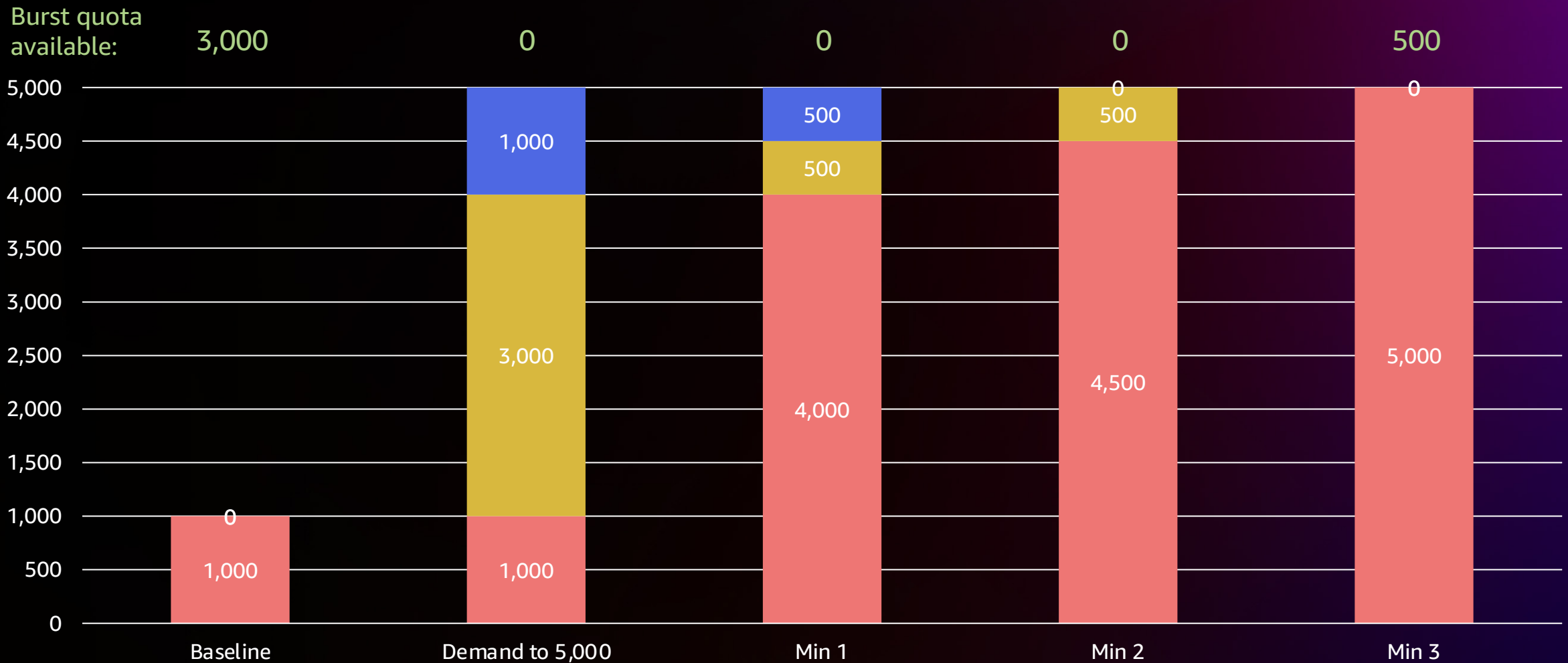
Concurrency limit = 1,000 (default, us-east-1)
Burst limit = 3,000 (default, us-east-1)

Burst quota
available:     3,000          3,000          3,000          3,000          3,000

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|

5,000

4,500

4,000

3,500

3,000              4,000          4,000          4,000          4,000

2,500

2,000

1,500

1,000        0              0              0              0              0
              1,000          1,000          1,000          1,000          1,000
500

0
          Baseline      Demand to 5,000      Min 1          Min 2          Min 3

■ Reqs Handled     ■ New Reqs Handled     ■ Reqs Throttled

Lambda bursting and ramp up

# What if you need more burst than the burst quota allows?

# Lambda concurrency controls

**Provisioned concurrency:**

- Sets floor on minimum number of execution environments

- Pre-warm execution environments to reduce cold-start impact

- Burst to use standard concurrency if desired
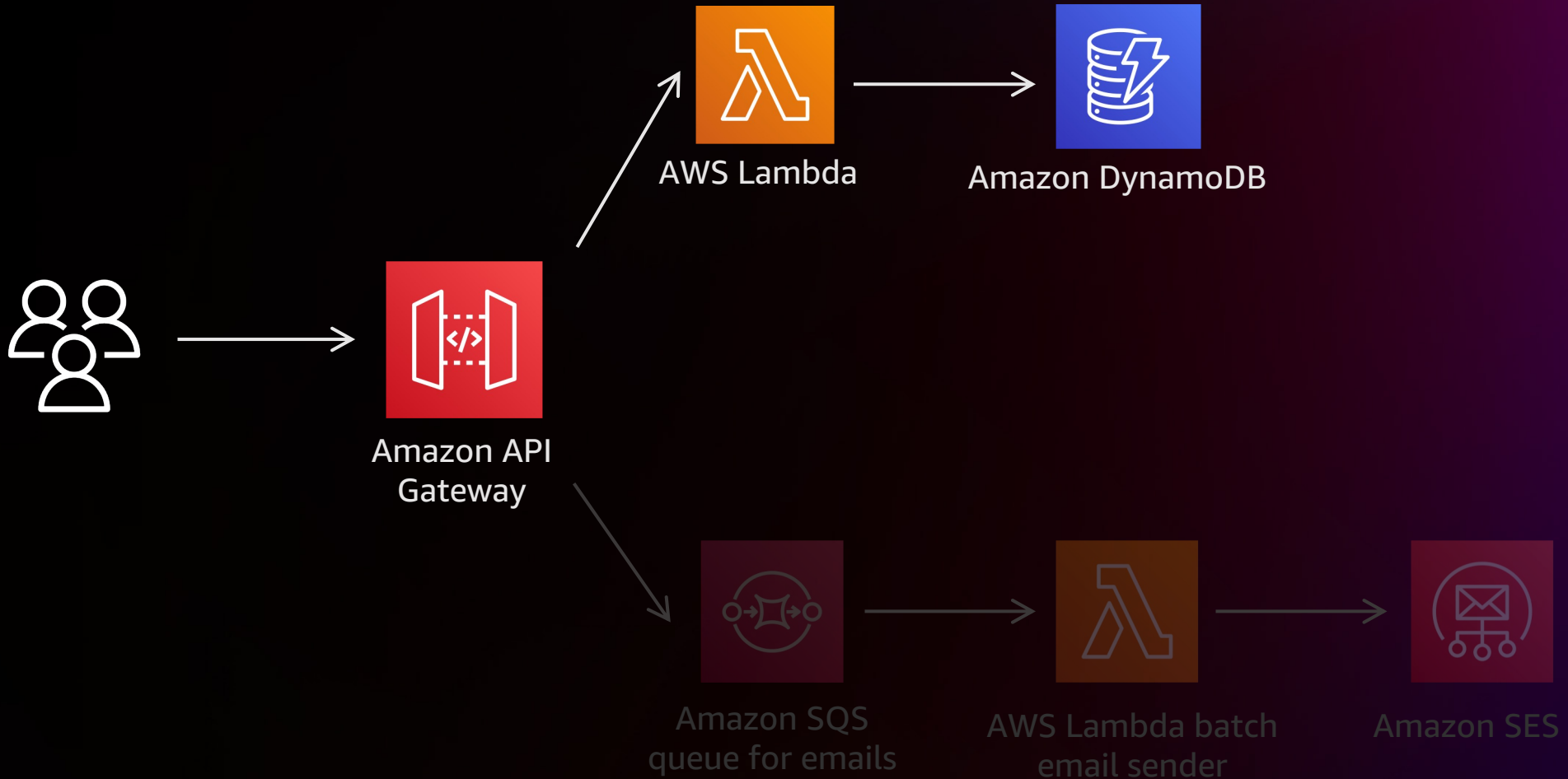
- Can save costs in certain situations



Function Scaling with Provisioned Concurrency

Requested provisioned concurrency

Burst limit

# Synchronous invocations

# Synchronous invocations: best practices

- **Load test**
  - Raise account-level concurrency quota in advance of bursts of demand
- **Set appropriate timeouts**
  - Set reasonable timeouts per Lambda function – do not set to 15 minutes
  - Configure API response timeout – default is 29 seconds
- **Retry with backoff**
  - Clients should retry, but with exponential backoff to not stampede the system
- **Implement idempotency**
  - Expect and account for retried operations
  - Application awareness and handling of scenarios is key

**Dive deeper at https://s12d.com/reInvent2020-Error-Handling**

# Asynchronous invocations

# Lambda function concurrency across models

**Amazon SNS**
*Async invocations*

No event store
"Fire and forget"
1 Lambda invocation per event

**Amazon SQS**
*Poller-based invocations*

Message delivery and processing durability
1 Lambda invocation per batch

# What if scaling too fast may overwhelm downstream systems?

# Lambda concurrency controls

## Reserved concurrency

- Sets ceiling on maximum number of execution environments – upper limit on maximum concurrency for a given function

- Also reserves that concurrency from the account's quota

**Reserved Concurrency**

Other functions

my-function-DEV

my-function-PROD

# Reserved concurrency

Account concurency



■ All functions

# Reserved concurrency



Account concurency

■ Other functions   ■ Function 1   ■ Function 2   ■ Function 3

**Reserved concurrency is an upper limit for a Lambda function's concurrency.**

**Reserved concurrency is always respected, regardless of integration.**

# How Lambda scales: Reserved Concurrency

*Time*

| 1 | Initialization | Execution |

| 2 | Initialization | Execution |

*Reserved Concurrency = 2*

# How Lambda scales: Reserved Concurrency



Time →

1 → Initialization | Execution

3 → Execution

2 → Initialization | Execution

3 → Throttle

Reserved Concurrency = 2

# Asynchronous invocations: best practices

- Invocation model
  - Leverage async invocations whenever possible
- Function configuration

  - Leverage batching for cost reduction and performance gains
  - Catch errors/failures and return altogether if using a batch size greater than 1
  - Configure Lambda on-failure destination to save and review failed messages
  - Monitor function error metrics to ensure they remain at or close to 0

- SQS queue configuration
  - Set the queue visibility timeout to at least 6x function timeout
  - Setup a DLQ with maxReceiveCount of at least 5

**SQS Lambda Best Practices:** https://s12d.com/SQS-Lambda-Best-Practices

# Concurrency tips

- Understand fundamentals

- Plan ahead and load test

- Monitor, monitor, monitor

  - `ConcurrentExecutions`

  - `Errors`

  - `Throttles`

  - SQS queue depth, Kinesis iterator age, etc.

# Open Discussion

# Check out these other sessions…

SVS404: A closer look at AWS Lambda

Wednesday (Nov 30) @ 10:00 AM

SVS401: Best practices for advanced serverless developers

Tuesday (Nov 29) @ 12:30 PM

SVS314: AWS Lambda performance tuning: Best practices and guidance

Thursday (Dec 1) @ 4:15 PM

# References

- Documentation

  - Invoke Lambda functions
    https://s12d.com/lambda-invocation

  - Using Lambda with Amazon SQS
    https://s12d.com/lambda-sqs

- Tools and guidance

  - AWS Lambda Powertools for Python
    https://s12d.com/powertools

  - Serverless Applications Lens - AWS Well-Architected Framework
    https://s12d.com/serverless-wa-lens

# Continue your AWS Serverless learning

## Learn at your own pace



Expand your Serverless skills with our Learning Plan on **AWS Skill Builder**

## Increase your knowledge



Use our **Ramp-Up Guides** to build your Serverless knowledge

## Earn AWS Serverless badge



aws
**Serverless**
2 0 2 2

Demonstrate your Knowledge by achieving **digital badges**



https://s12d.com/serverless-learning

# Serverless Land for more resources



**Session slides and links:** https://s12d.com/reInvent2022-Lambda-Concurrency

# Please complete the session survey

# Thank you!

Brian Zambrano (he/him)

🐦 @brianzambrano

Justin Pirtle (he/him)

🐦 @justinpirtle

# The function lifecycle

Download **your code**   Start new **execution environment**   Bootstrap the **runtime**   Start your **code**

| Full cold start | Partial cold start | Warm start |

AWS **optimization**   Your **optimization**

# CloudWatch metrics: Lambda

What to monitor during testing and in production:

- *Errors* - Number of failed function invocations due to timeout, code exceptions, or service issues

- *Throttles* - Number of throttled function invocations caused by account/burst concurrency limit or intentional reserved concurrency ceiling

- *ConcurrentExecutions* – Count of function execution environments simultaneously processing events. View as Maximum for give time window

- *ProvisionedConcurrentExecutions* – Count of function execution environments simultaneously processing events on provisioned concurrently

# Lambda concurrency quota reached

# CloudWatch metrics: Lambda errors/throttles



Lambda 429 TooManyRequests throttles surfaced as:
- 500 status code response from REST API Gateway deployment
- 503 status code response from HTTP API Gateway deployment

# Test: Effect of execution duration on concurrency

- Send 100 RPS to APIGW

- No delay between requests

- Change execution duration every 5 mins via query param

  - 5ms
  - 50ms
  - 500ms
  - 1000ms

# Test: Effect of RC with various RPS

- Set reserved concurrency to 50

- Average request duration 1000ms (1s)

- Send various TPS to APIGW

  - 25 RPS

  - 50 RPS

  - 75 RPS

- No delay between requests

# Test: Synchronous Lambda invocations

- Account/region quota: 10,000 (us-west-2)

- Direct Lambda sync invoke with Go

- Launch three batches with various concurrency every 5 mins
  - t1: 1000
  - t1 + 5m: 4000 (5000 total)
  - t1 + 10m: 5000 (10000 total)

- No delay between requests

# Test: Lambda scaling with long SQS queue

- Account/region concurrency quota: 10,000 (us-west-2)

- Put an artificial delay in Lambda function

- Load a lot of messages onto standard queue

# Test: Lambda with reserved concurrency and SQS

- Load 3,000 messages onto standard queue

- Use various reserved concurrency settings:

  - 200

  - 60

  - 30

  - 15

- Throttling begins at RC=60
- Lower RC == more throttling
- 100% of messages processed

# Test: Lambda with reserved concurrency and SQS

- Load 10,000 messages onto standard queue

- Set reserved concurrency to 15

- Insert longer artificial delay

100% of messages processed

Setup a DLQ with the recommended settings:

maxReceiveCount >= 5

VisibilityTimeout = 6 * function timeout

# Amazon SQS to Lambda: Scaling Invocations

# Amazon SQS to Lambda: Scaling Invocations

Send message

Long poll

D    C    B    A

Lambda

Poller

# Amazon SQS to Lambda: Scaling Invocations

Send message

Long poll    Batch of 3

Lambda

D    C    B    A

Poller

# Amazon SQS to Lambda: Scaling Invocations

# Amazon SQS to Lambda: Scaling Invocations



Send message

Lambda

1 invocation per received batch

Poller

Function instance

Batch of 3

D  C  B  A

# Amazon SQS to Lambda: Scaling Invocations



Send message

Batch of 1

Lambda

Poller

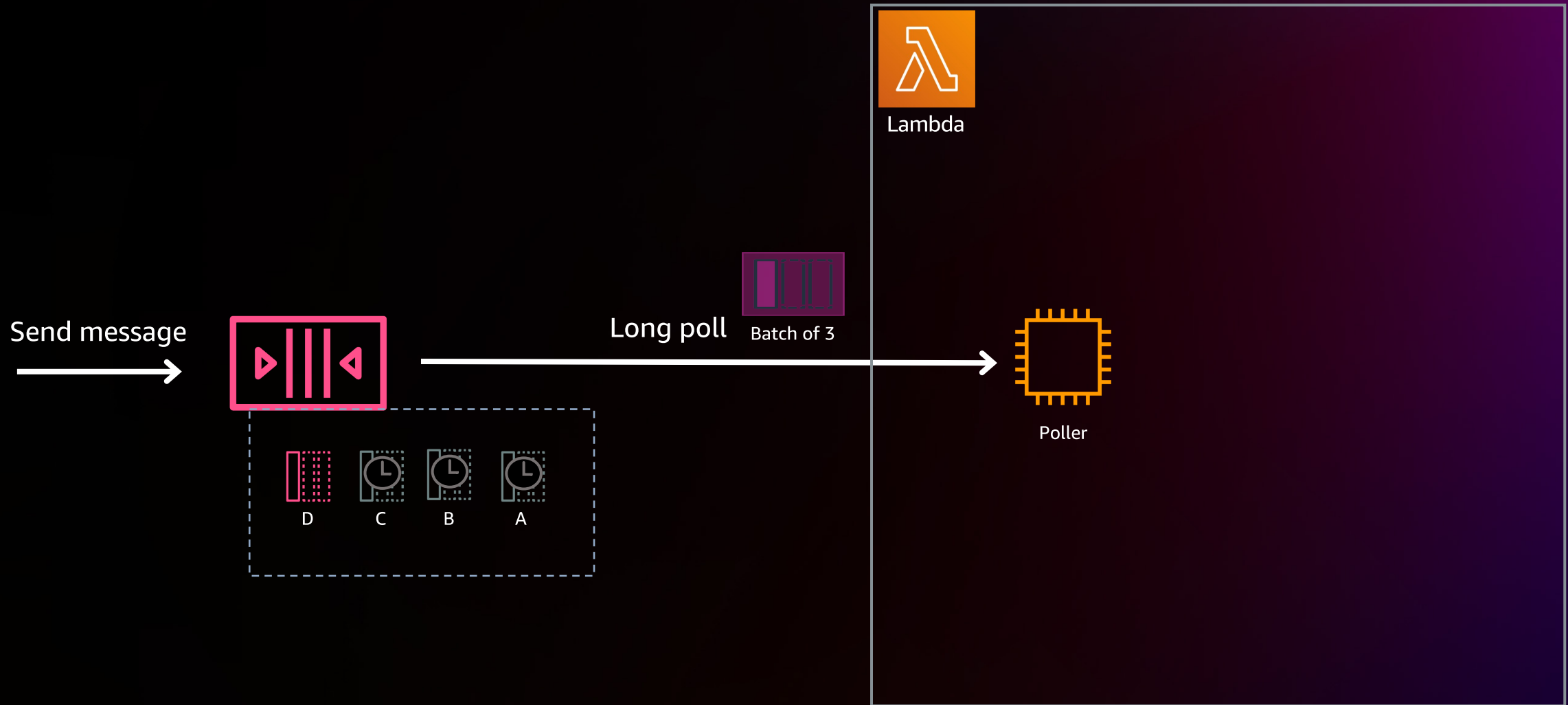Function instance

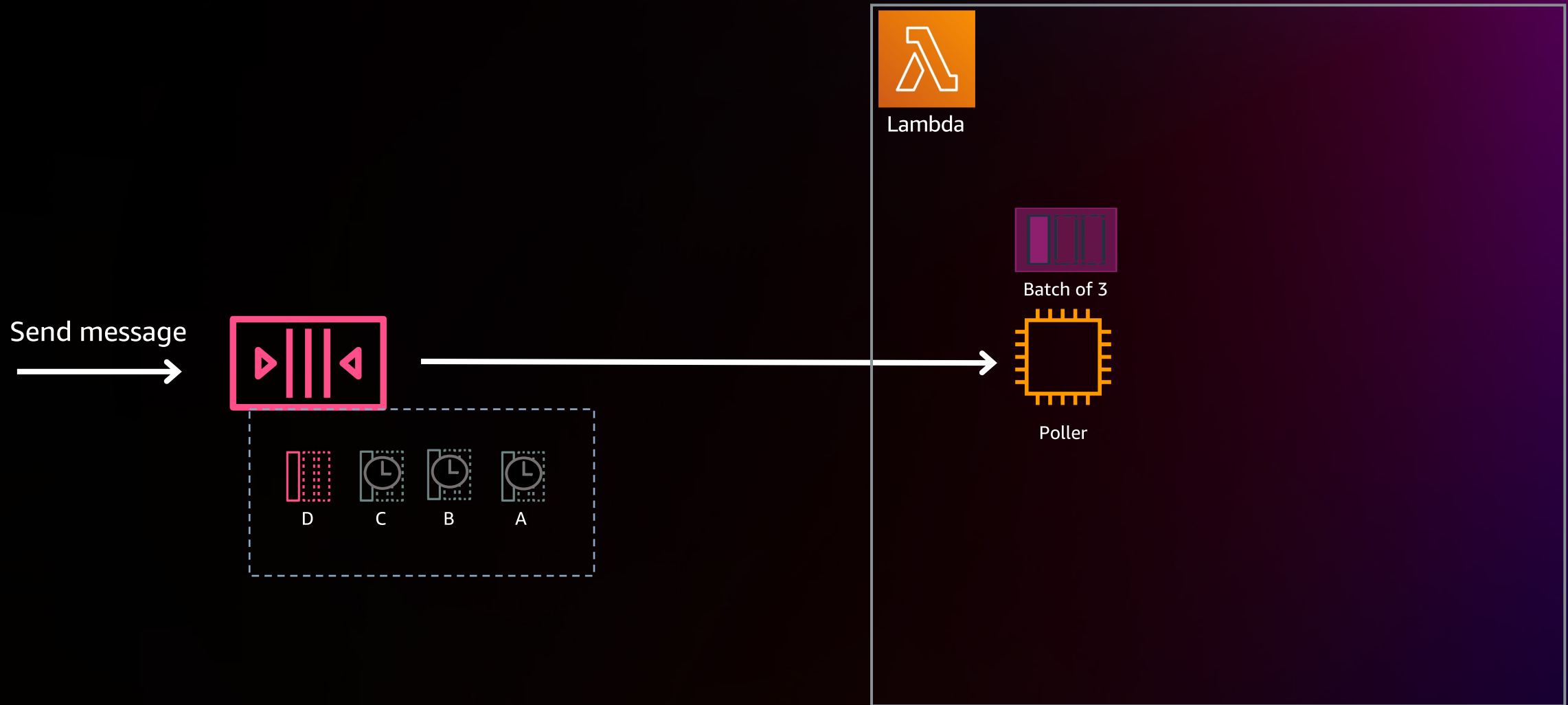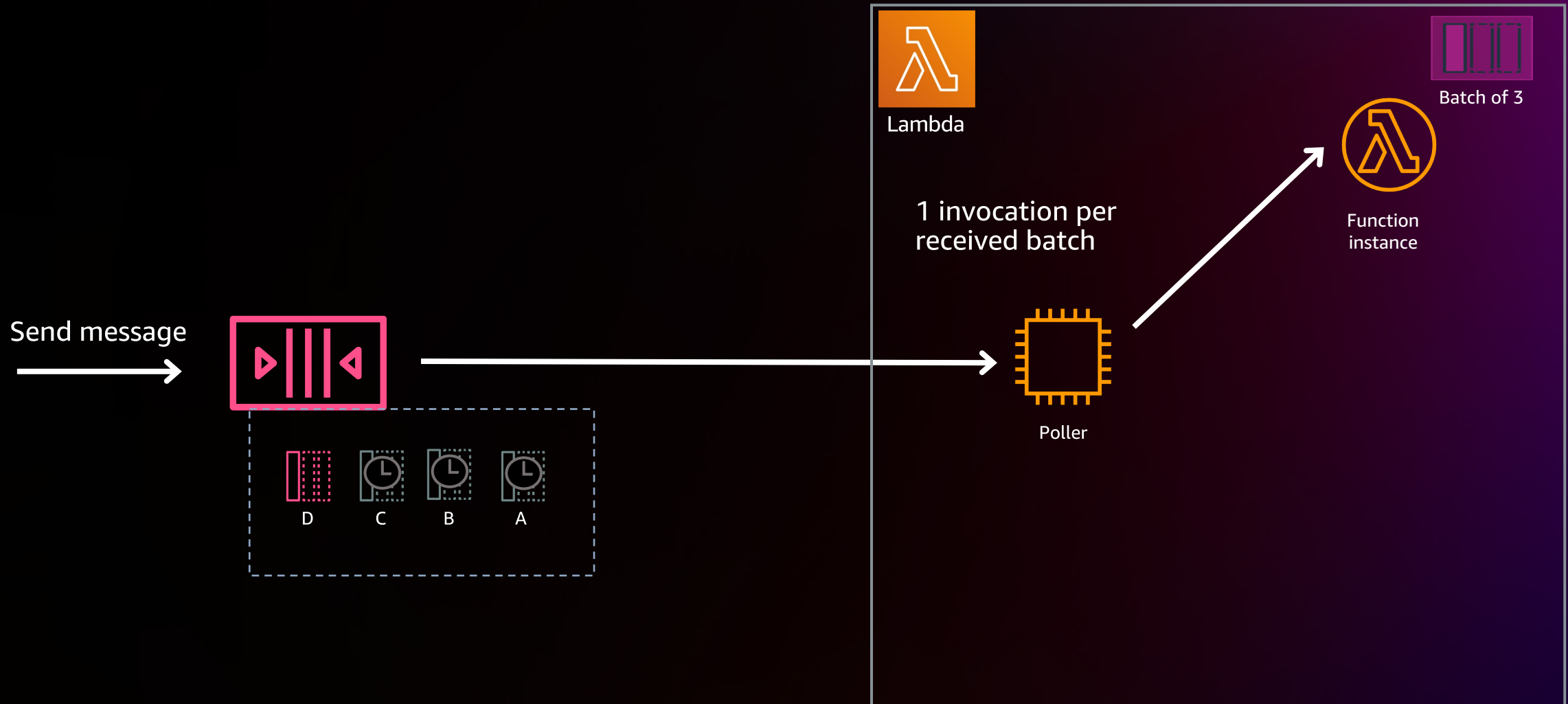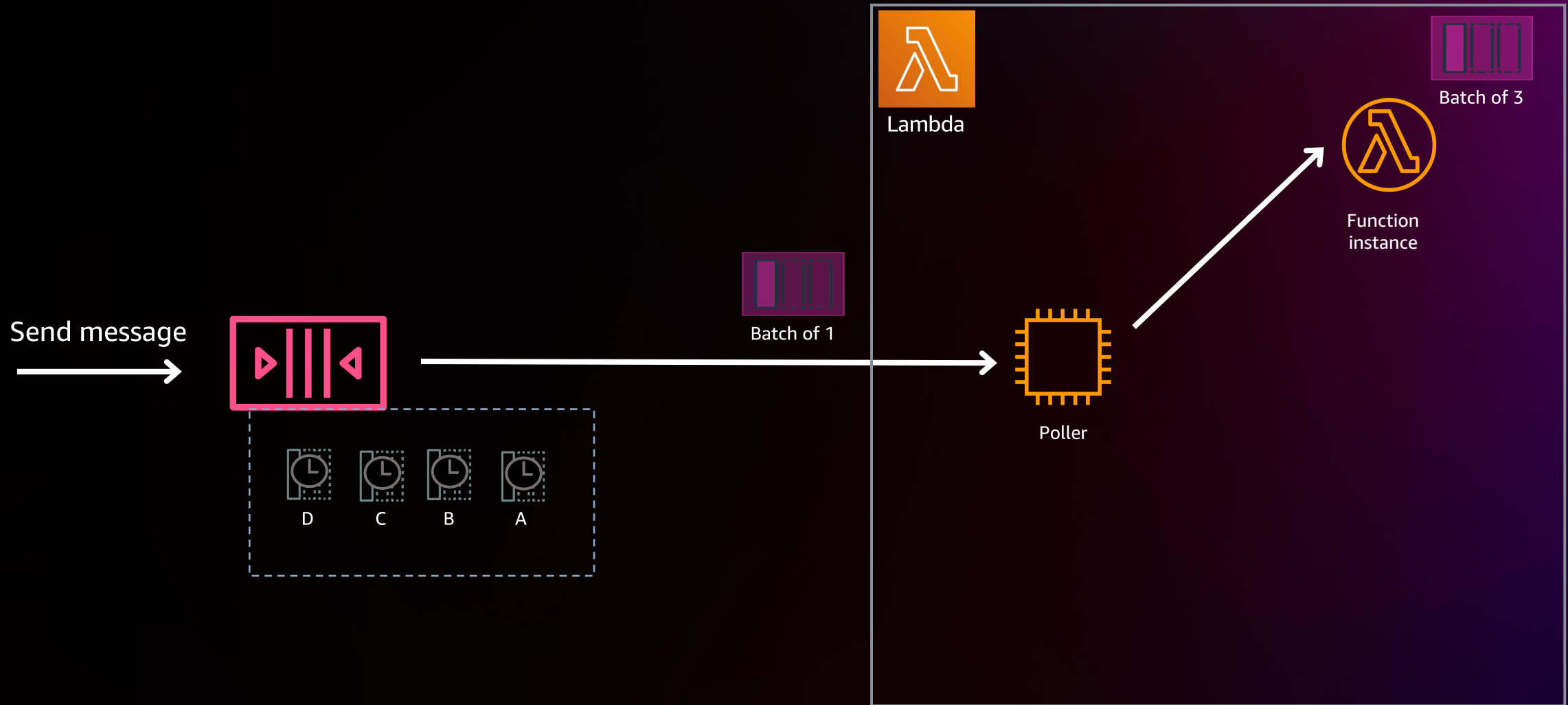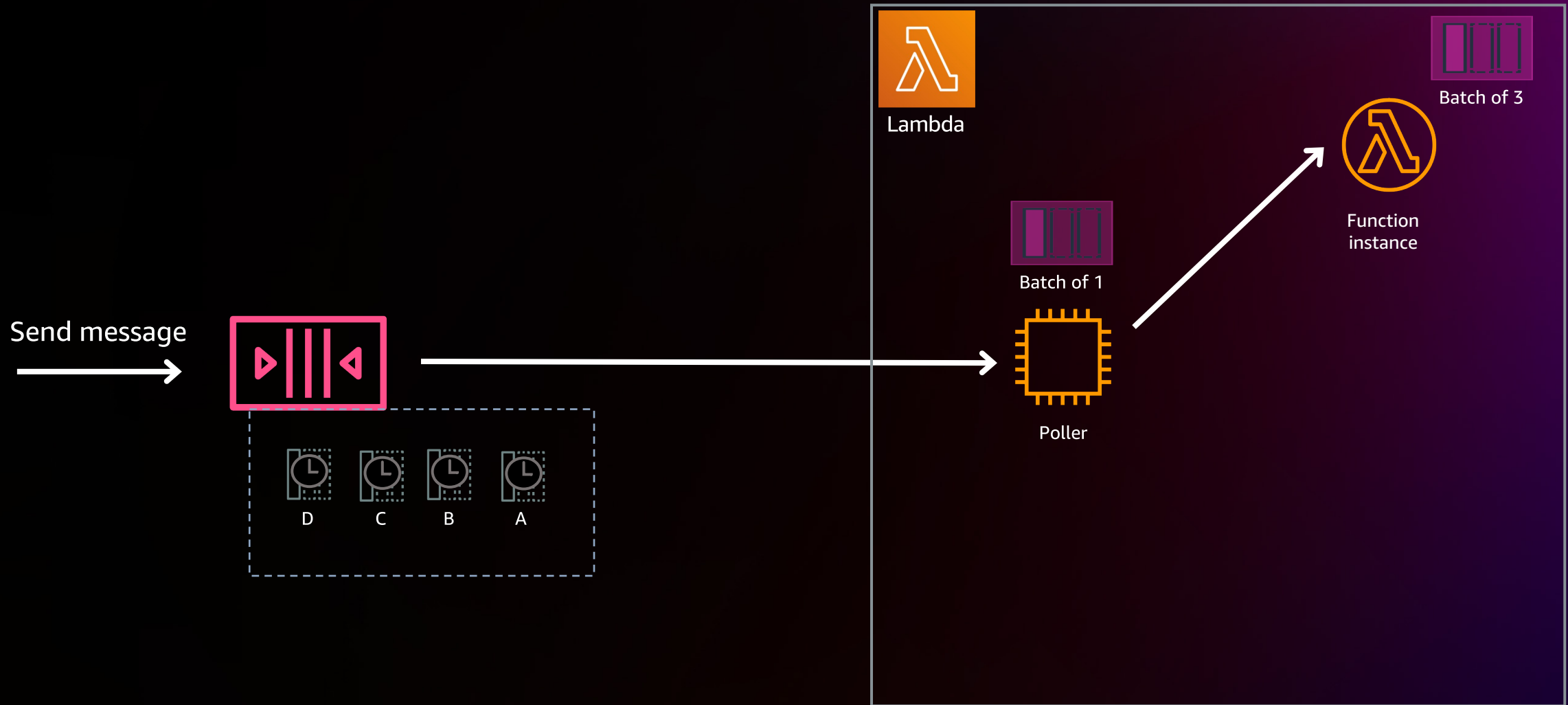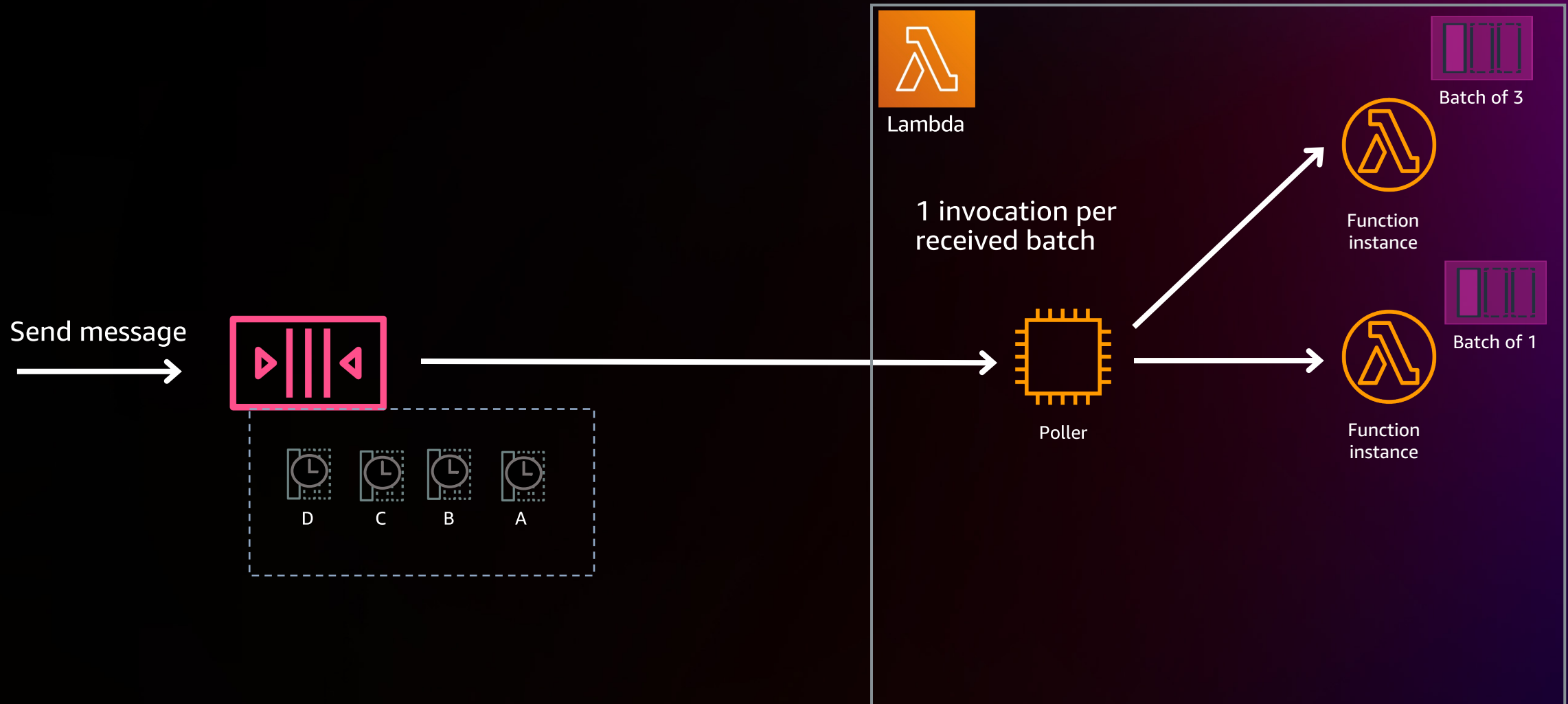Batch of 3

D    C    B    A

# Amazon SQS to Lambda: Scaling Invocations

# Amazon SQS to Lambda: Scaling Invocations

# Amazon SQS to Lambda: Scaling Invocations

Send message

Lambda

D   C   B   A

Poller

Function instance

Batch of 3

Function instance

Batch of 1

# Amazon SQS to Lambda: Scaling Invocations

Send message

Lambda

Batch of 3

Function instance

Poller

Batch of 1

Function instance

D  C  B  A

# Amazon SQS to Lambda: Scaling Invocations

# Amazon SQS to Lambda: Scaling Invocations



Lambda

Send message

Poller

Function
instance

Batch of 1

Function
instance

D

# Amazon SQS to Lambda: Scaling Invocations

Send message

Lambda

Poller

Function instance

Function instance

D

# Amazon SQS to Lambda: Scaling Invocations



Send message

Lambda

Poller

Function instance

Function instance

D

# Amazon SQS to Lambda: Scaling Invocations



Lambda

Send message

Batch of 1

Receive count = 2

D

Poller

Function instance

Function instance

# Kinesis Data Streams



Up to five Lambda functions can be subscribed to a single stream*

Lambda service polls the stream for new records once per second

Data producer

Kinesis Data Streams

Lambda service

Lambda function A

Lambda function B

Data producers continually add new records to the stream

When new records are available, the subscribed Lambda functions are invoked with new records

*Kinesis Data Analytics and Kinesis Data Firehose also count toward this subscriber limit

# DynamoDB Streams

Up to two Lambda functions can be subscribed to a single DynamoDB stream
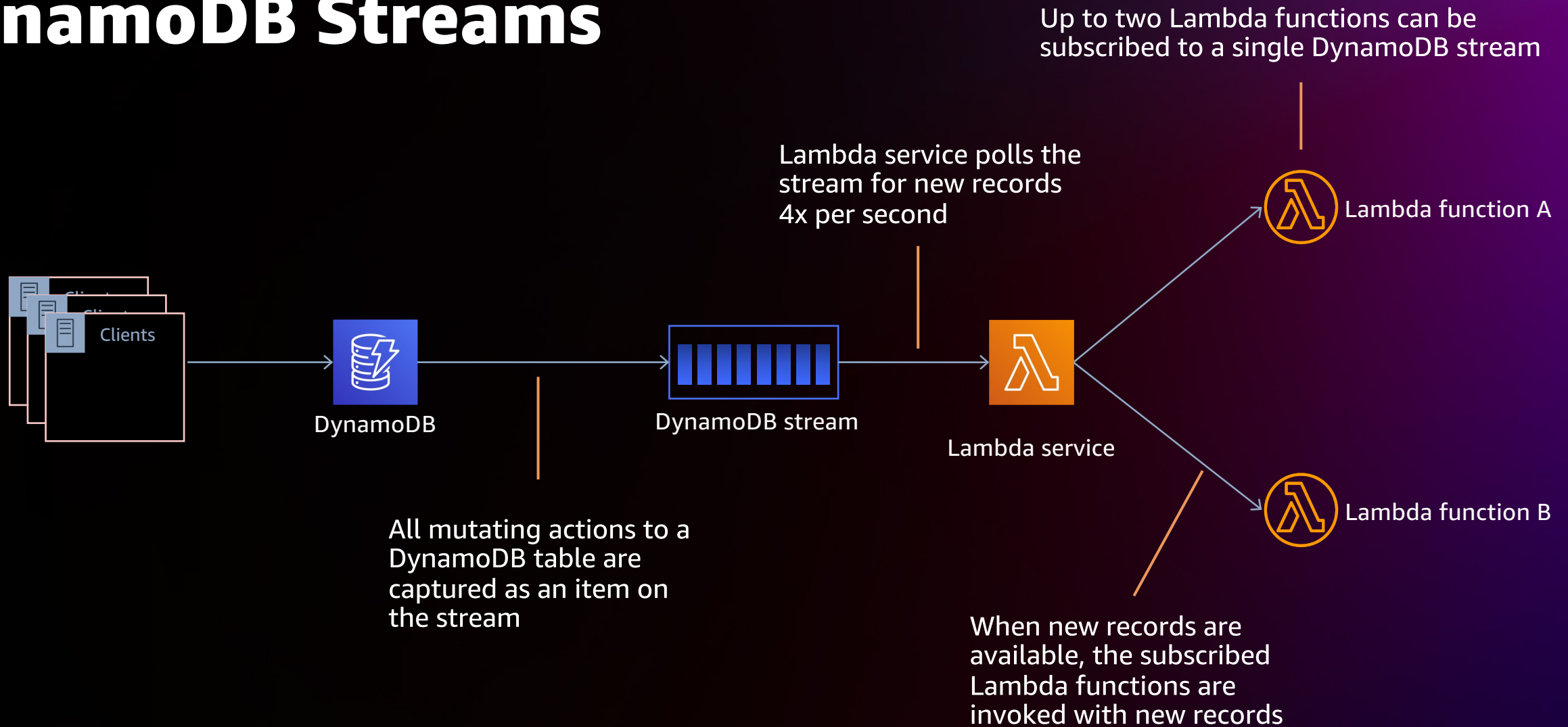
Lambda service polls the stream for new records 4x per second

Clients

DynamoDB

DynamoDB stream

Lambda service

Lambda function A

Lambda function B

All mutating actions to a DynamoDB table are captured as an item on the stream

When new records are available, the subscribed Lambda functions are invoked with new records
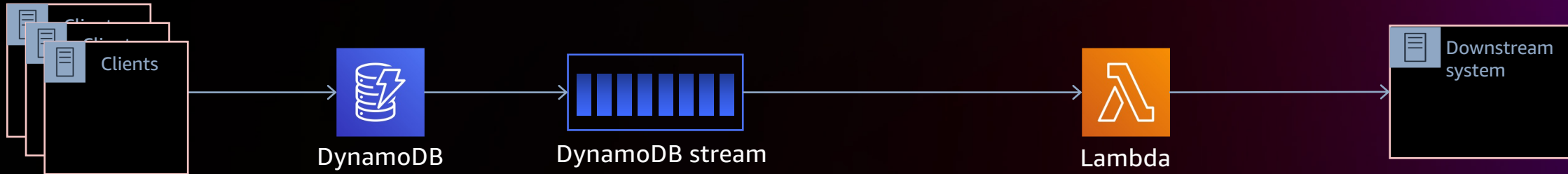
# Kinesis data stream monitoring



## Kinesis Metrics/Alarms

- GetRecords.IteratorAgeMilliseconds

- IncomingRecords/IncomingBytes

- ReadProvisionedThroughputExceeded

- WriteProvisionedThroughputExceeded

## Lambda Metrics/Alarms

- Errors

- IteratorAge

- Throttles

# DynamoDB stream monitoring



Clients

DynamoDB

DynamoDB stream

Lambda

Downstream system

## DynamoDB Streams metrics/alarms

- ReturnedRecordsCount/ReturnedBytes

- UserErrors

## Lambda metrics/alarms

- Errors

- IteratorAge

- Throttles

- Duration