AWS
re:Invent
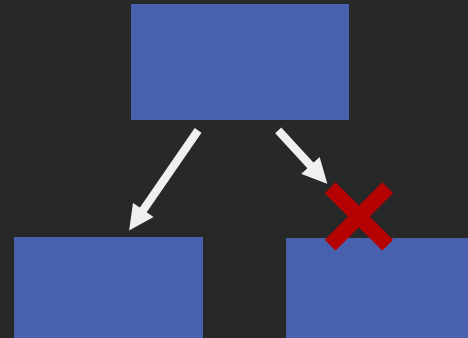
# Table of contents



**Load shedding**

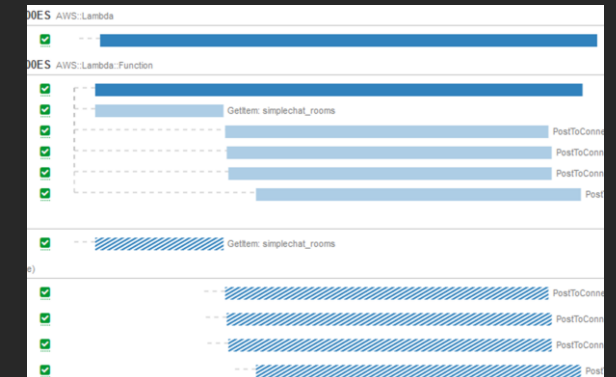Avoid brownout by rejecting excess load



**Dependency isolation**

Prevent one dependency from affecting unrelated functionality



**Avoiding queue backlogs**

Prevent a backlog from extending recovery time



**Operating**

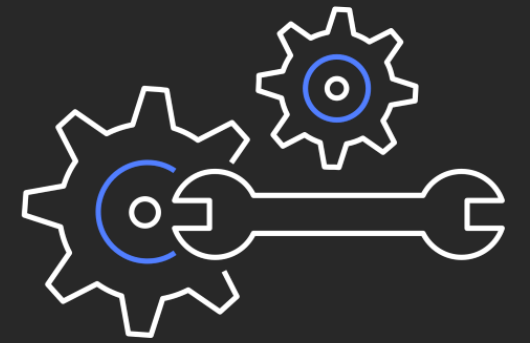Quickly diagnose and mitigate issues

# Structure

Theory

Serverless
resiliency tools

Motivation for
serverless

AWS under
the hood

# Overload, illustrated

"Life imitates art far more than art imitates life"

- Oscar Wilde, 1889

algorithms

algorithms

"Life imitates ~~art~~ far more than ~~art~~ imitates life"

algorithms

algorithms

"Life imitates ~~art~~ far more than ~~art~~ imitates life"

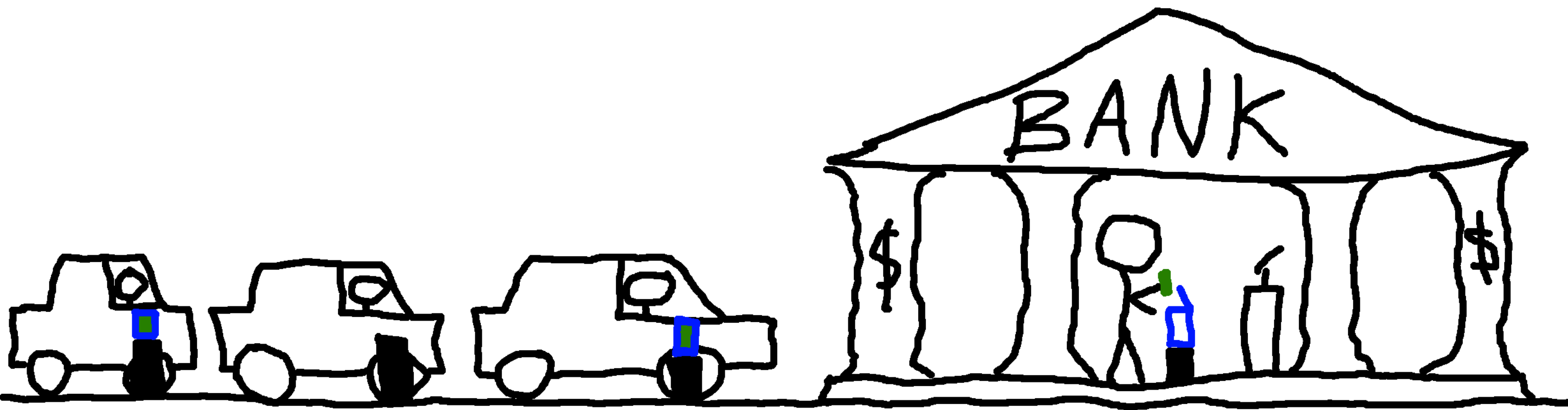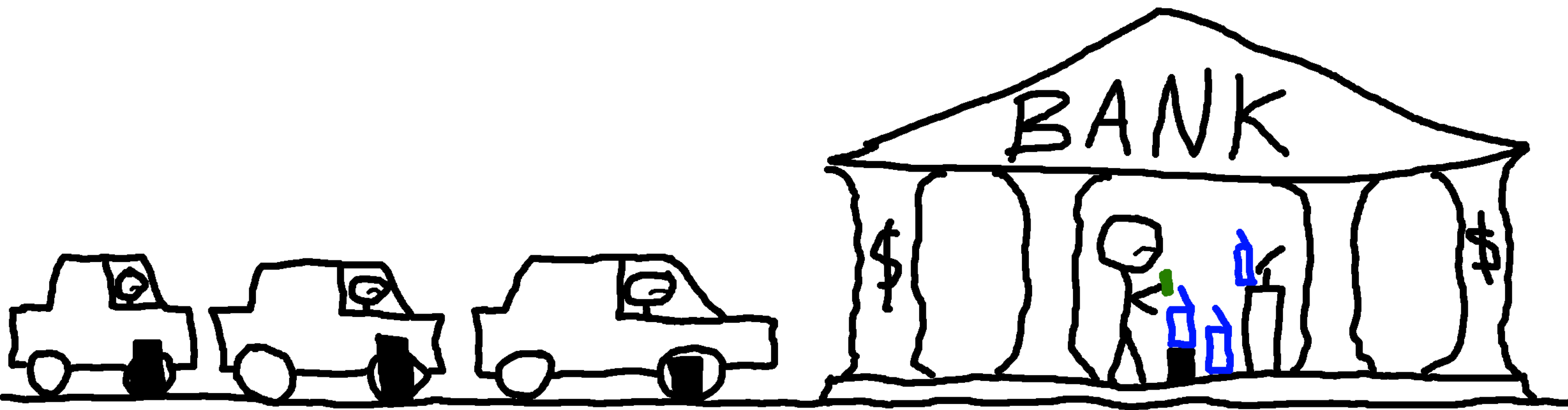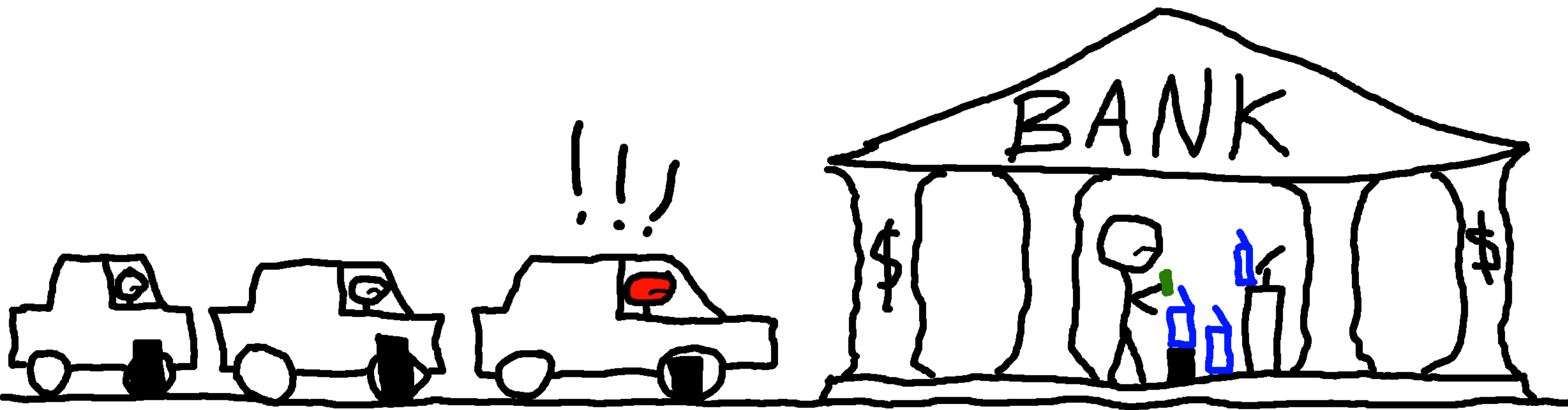"Algorithms imitate life far more than life imitates algorithms"

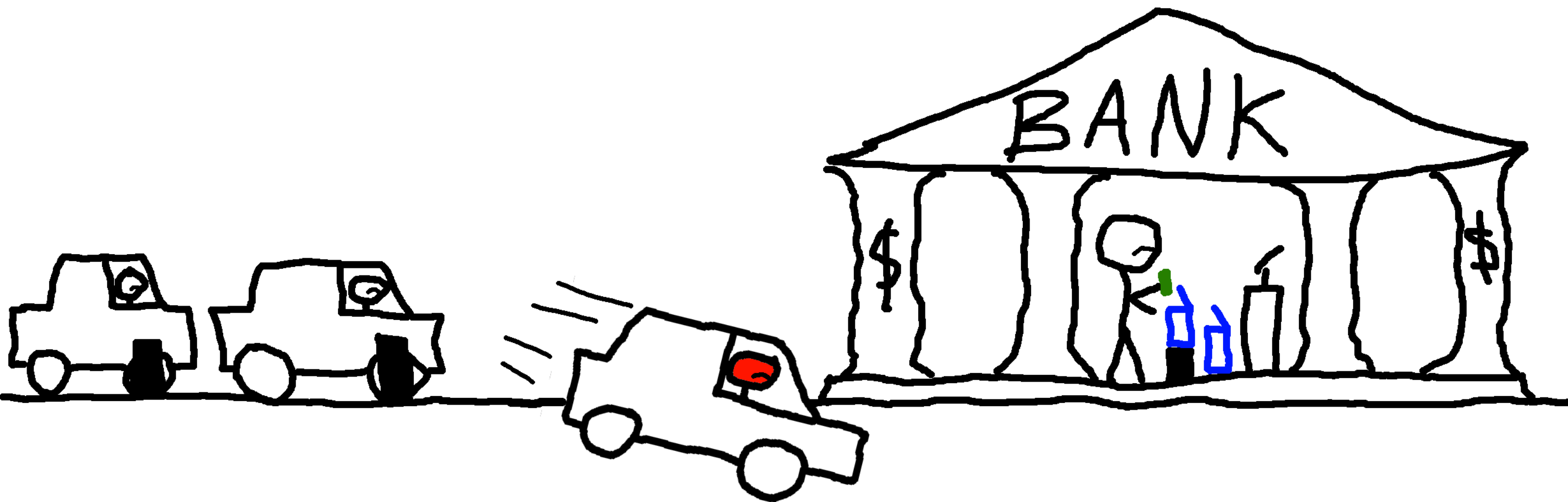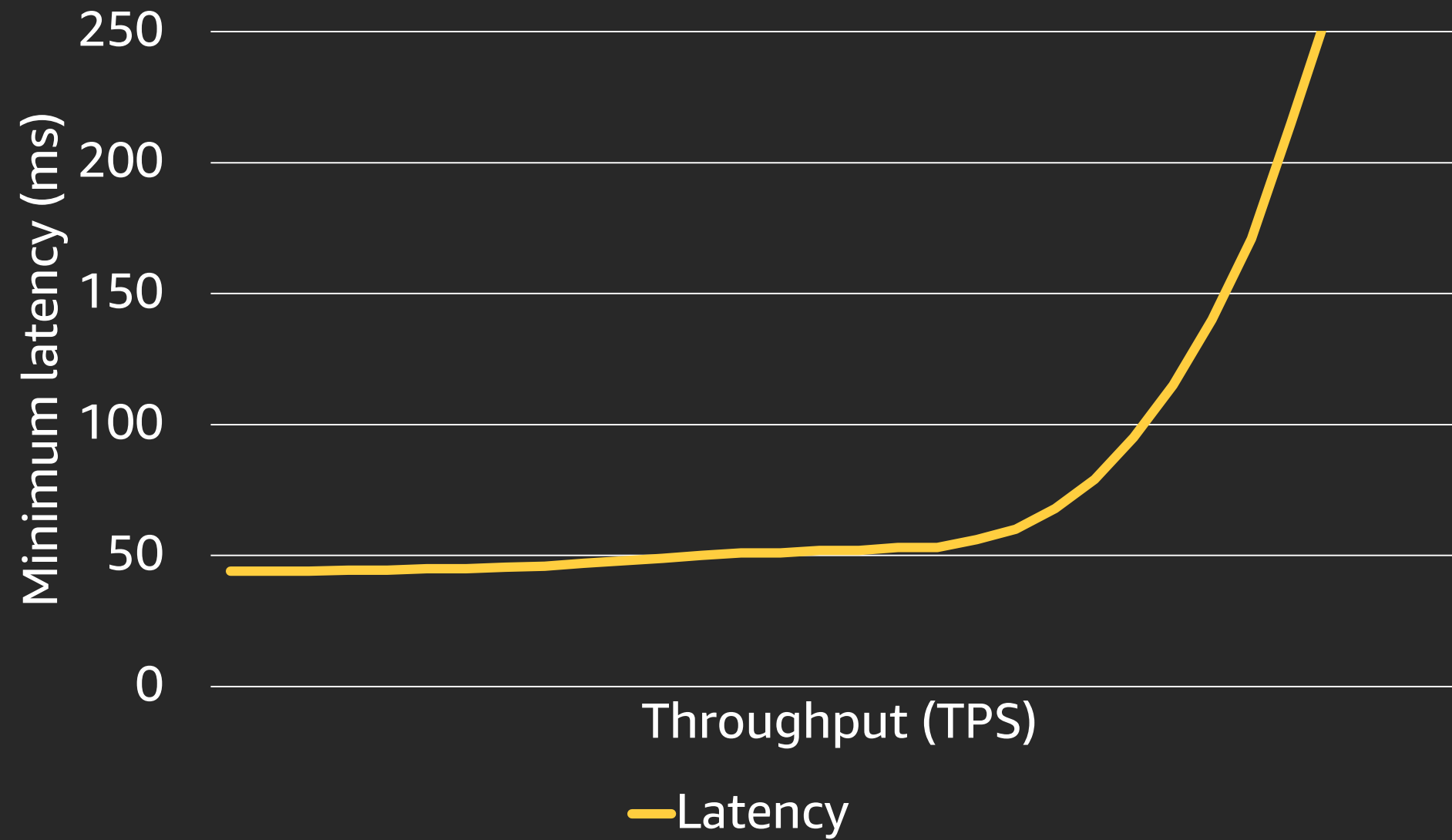– Somebody, probably
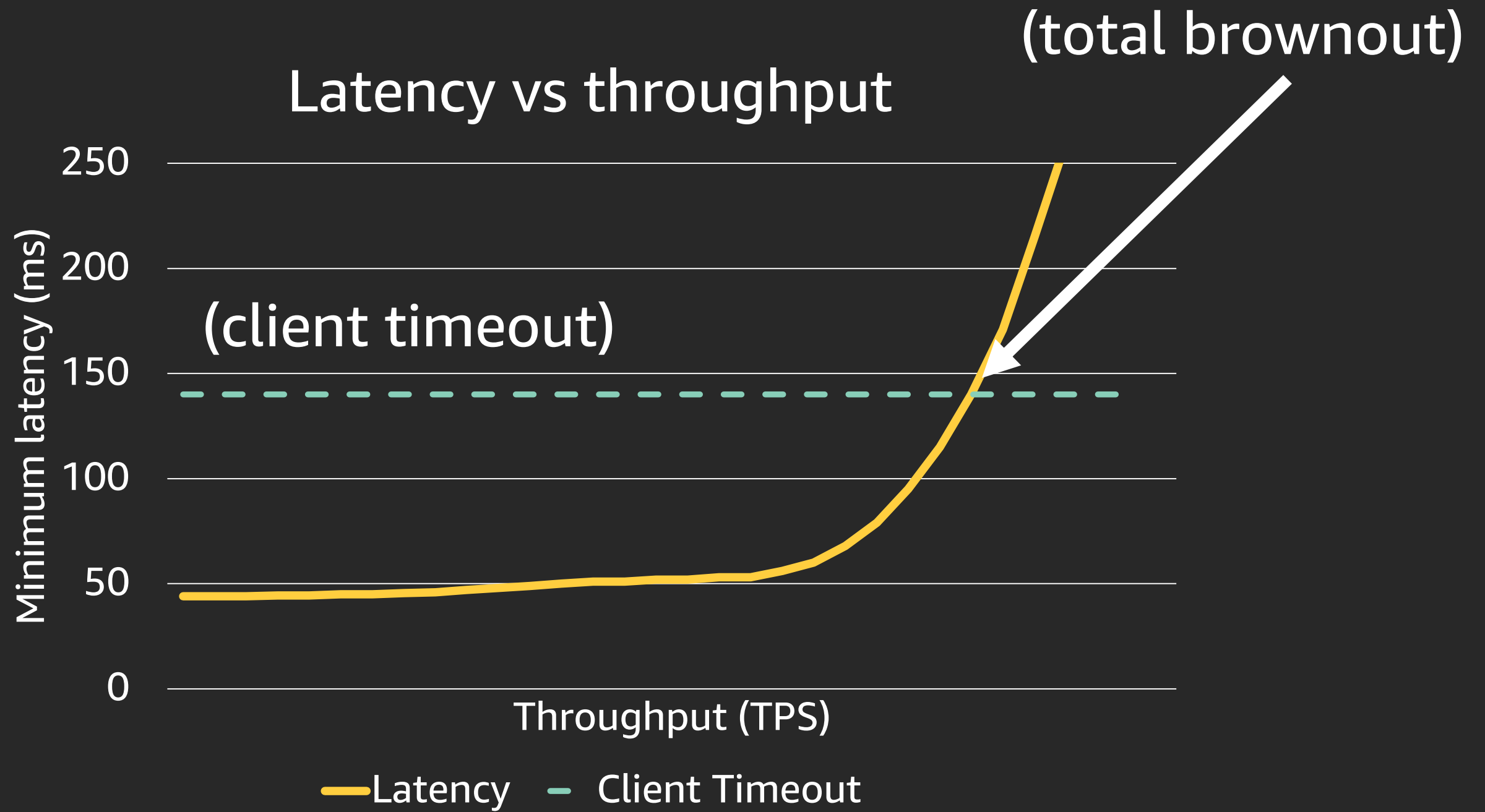
# Back in high school…

# Amdahl's Law, Universal Scalability Law

You can parallelize a system up to the point where contention becomes the bottleneck

# Latency vs throughput



**(total brownout)**

**(client timeout)**

250

200

150

100

50

0

Minimum latency (ms)
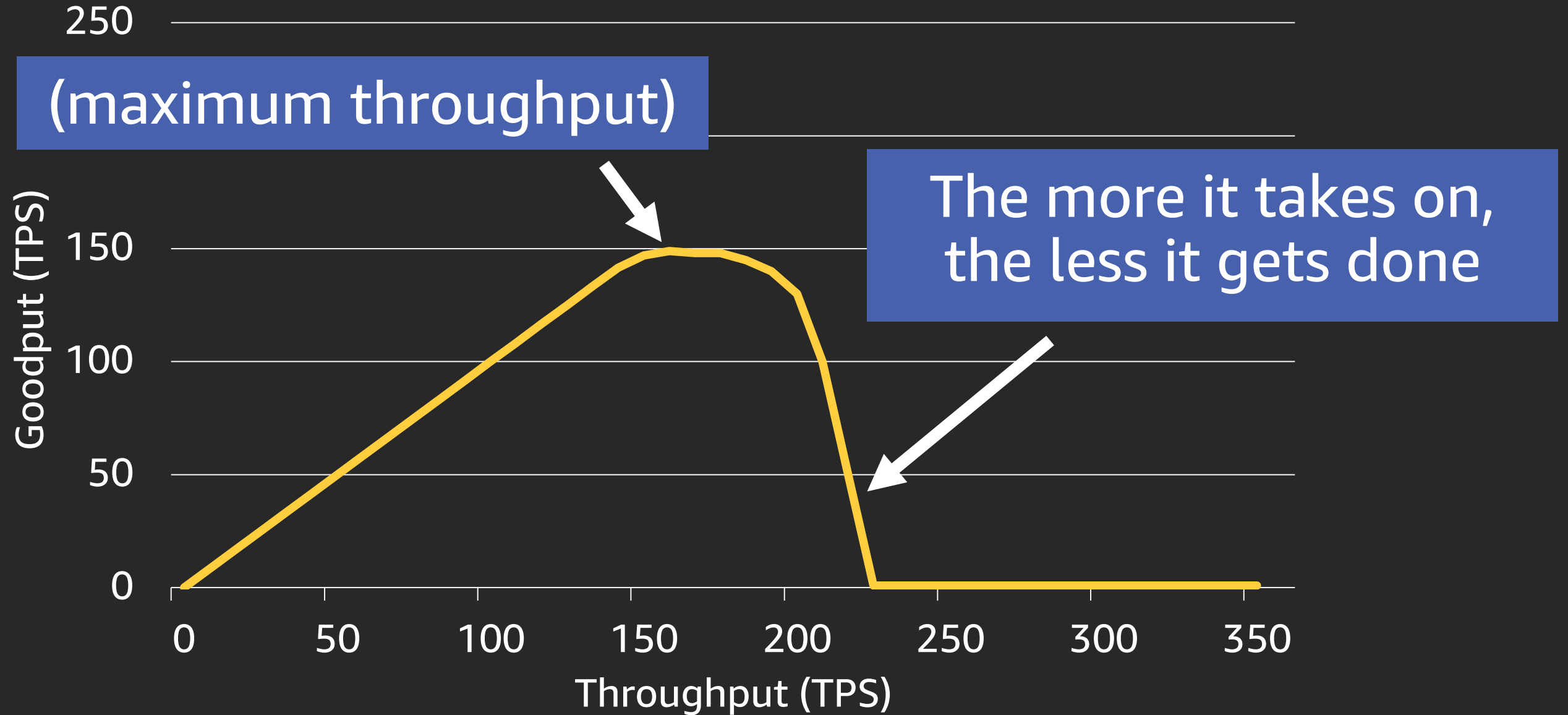
Throughput (TPS)

— Latency　－ － Client Timeout
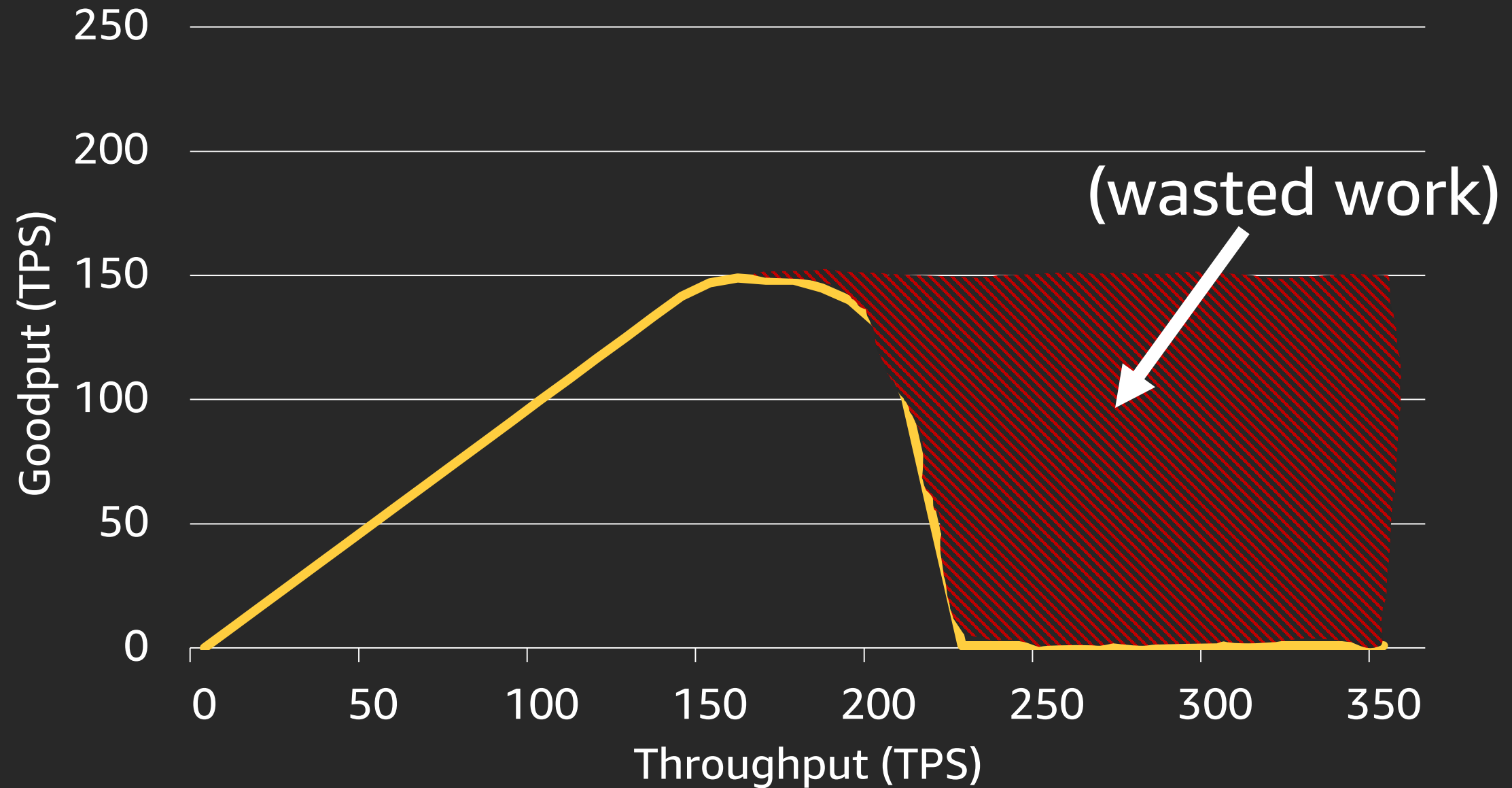
# Goodput vs throughput

Goodput vs throughput
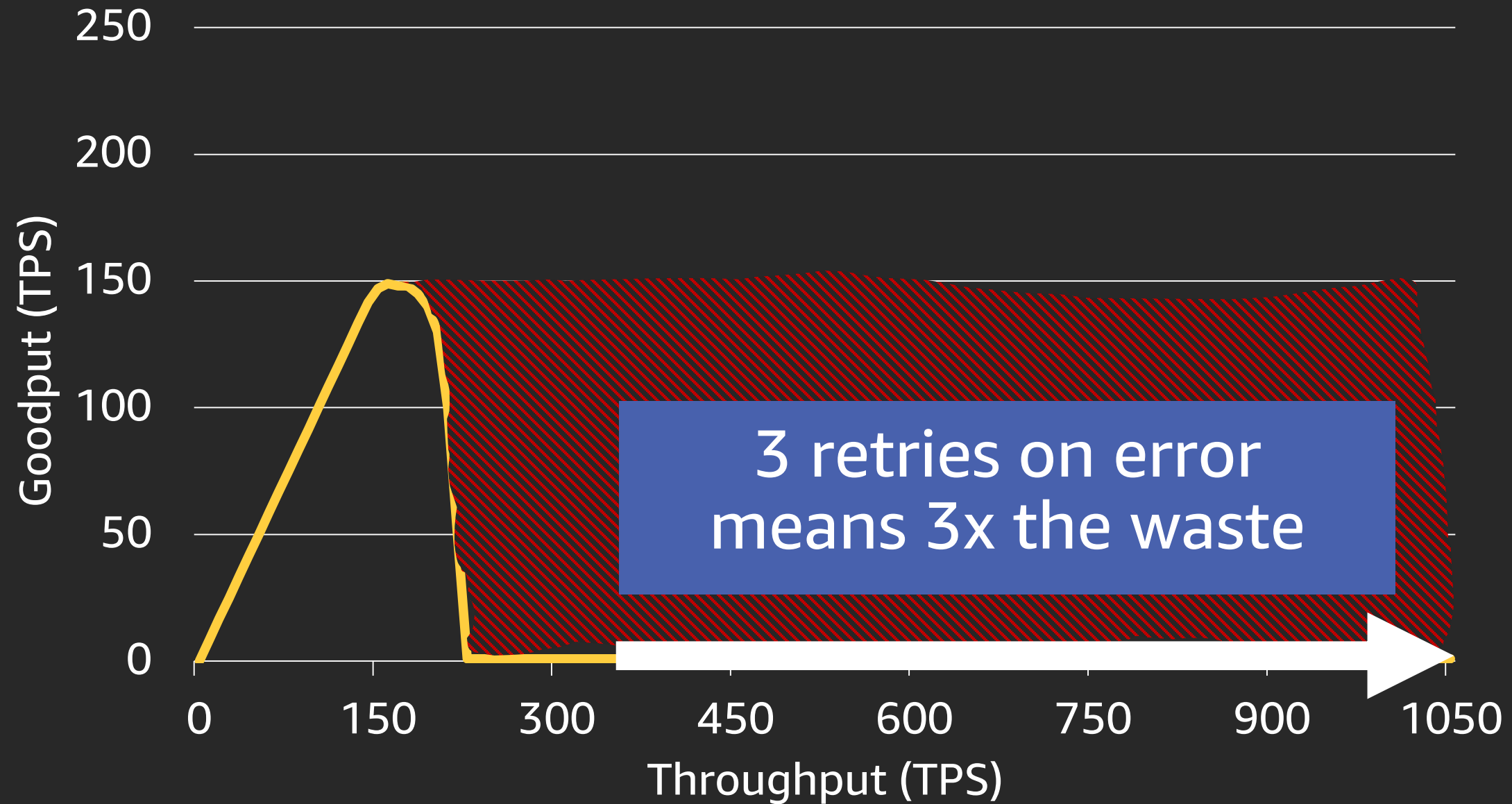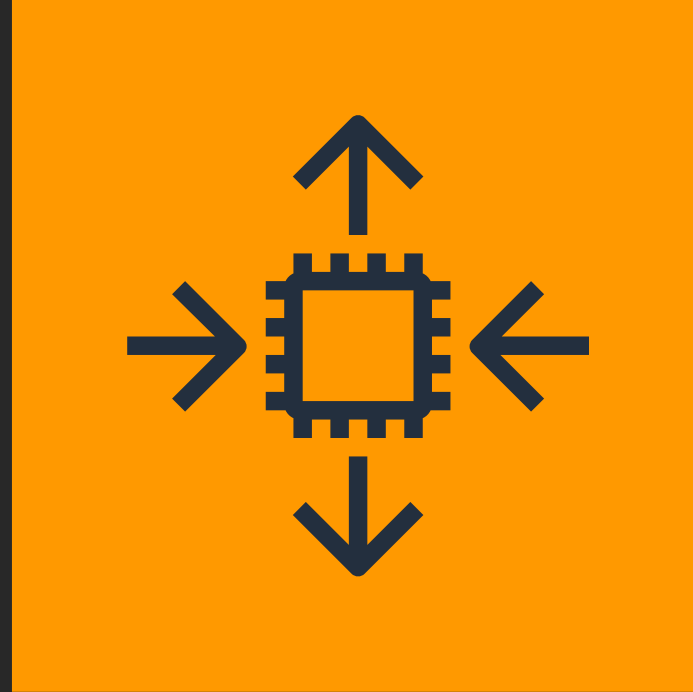
Servers are too optimistic

Goodput vs throughput
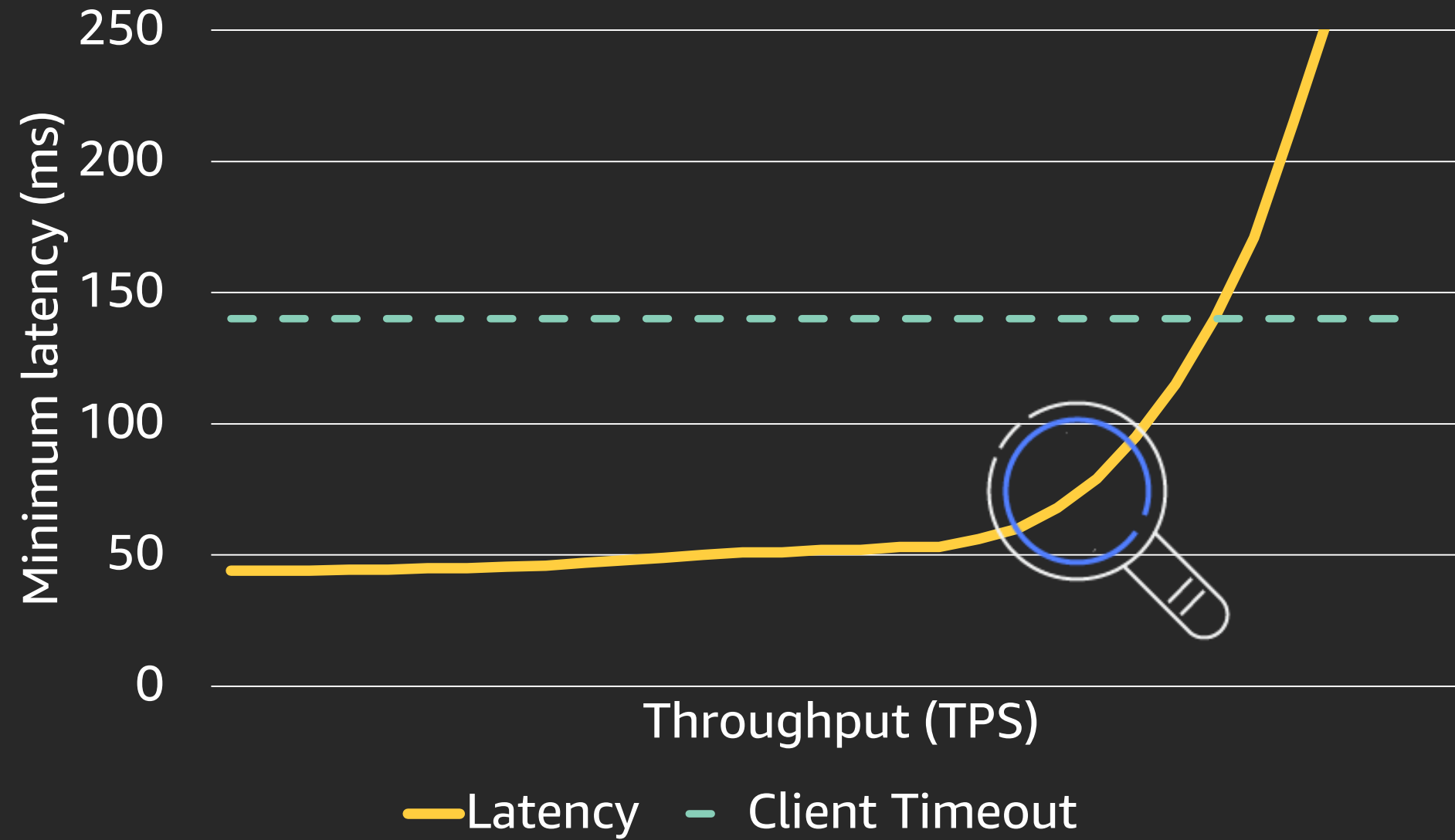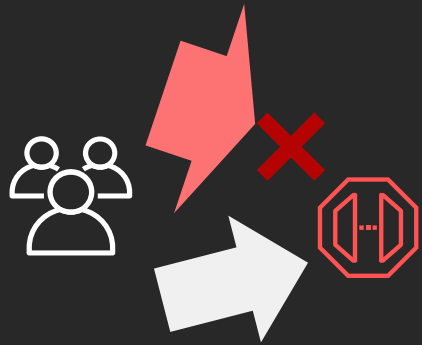
Clients retry

Goodput vs throughput

Goodput (TPS)

Throughput (TPS)

3 retries on error means 3x the waste

Amazon EC2
Auto Scaling

# Understand tipping points

# Latency vs throughput



Minimum latency (ms)

250

200

150

100

50

0

Throughput (TPS)
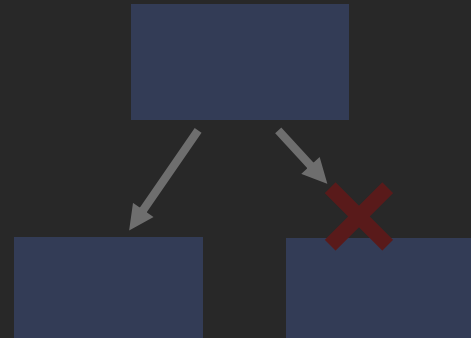
— Latency  – – Client Timeout
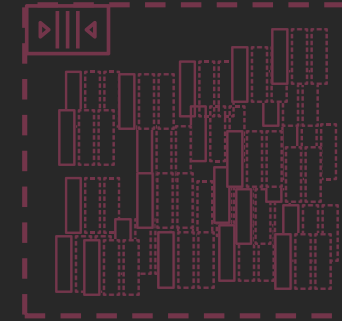
# Chapter one

**Load shedding**

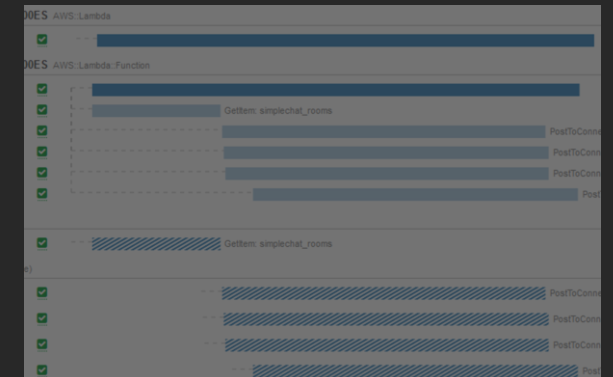**Avoid brownout by rejecting excess load**

Dependency isolation

Prevent one dependency from affecting unrelated functionality

Avoiding queue backlogs
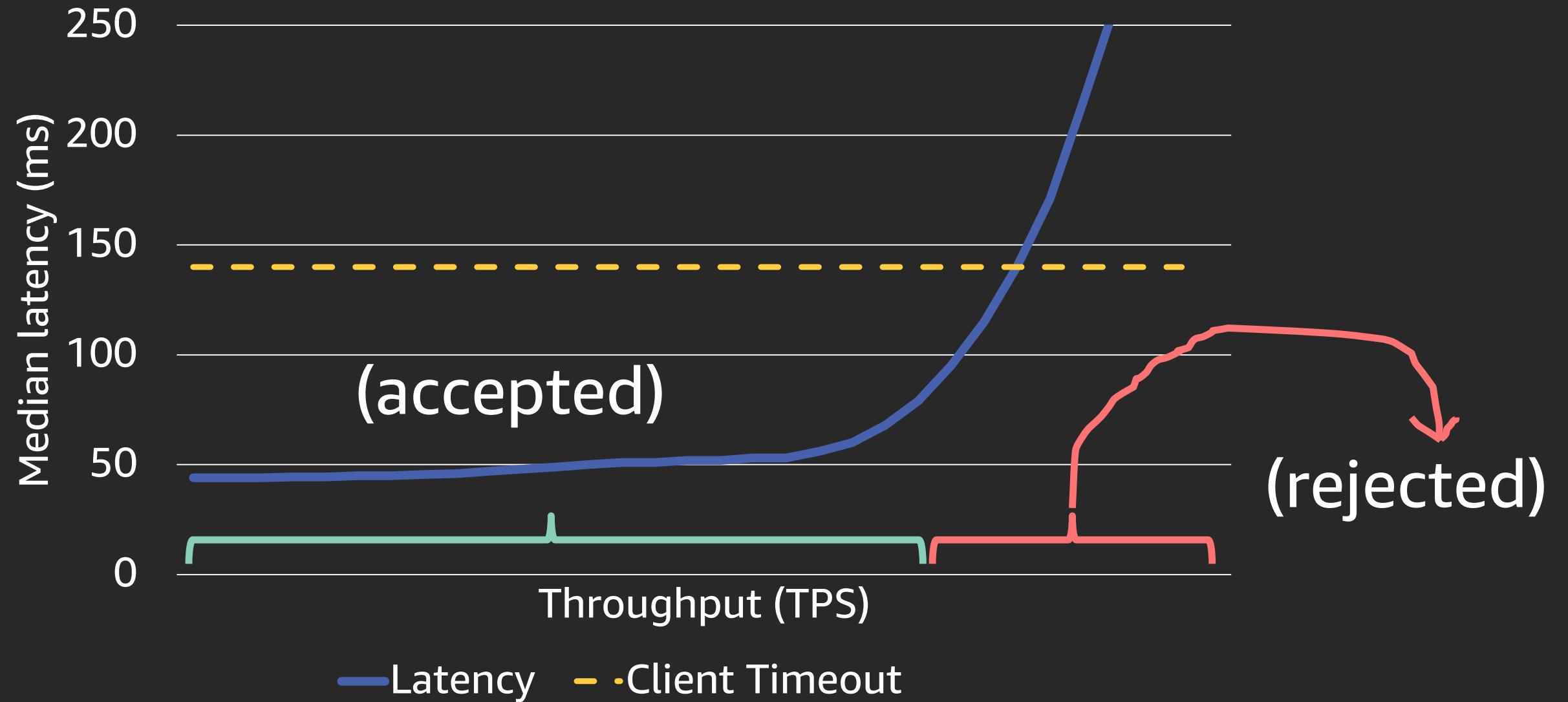
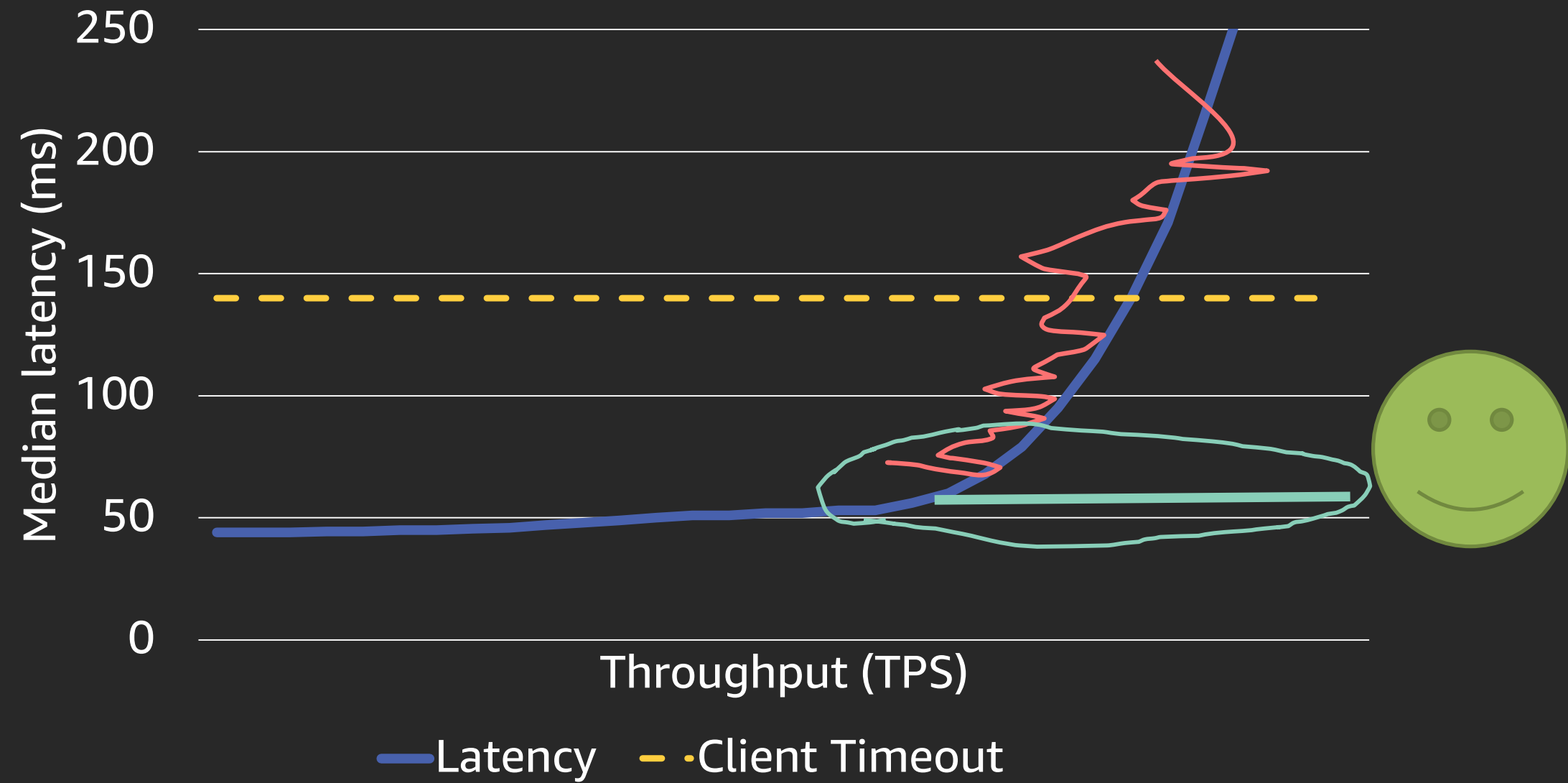Prevent a backlog from extending recovery time

Operating

Quickly diagnose and mitigate issues

Cheaply reject excess work

Latency vs throughput

(accepted)

(rejected)

Median latency (ms)

Throughput (TPS)

Latency    Client Timeout

# Latency vs throughput

**Median latency (ms)** vs **Throughput (TPS)**

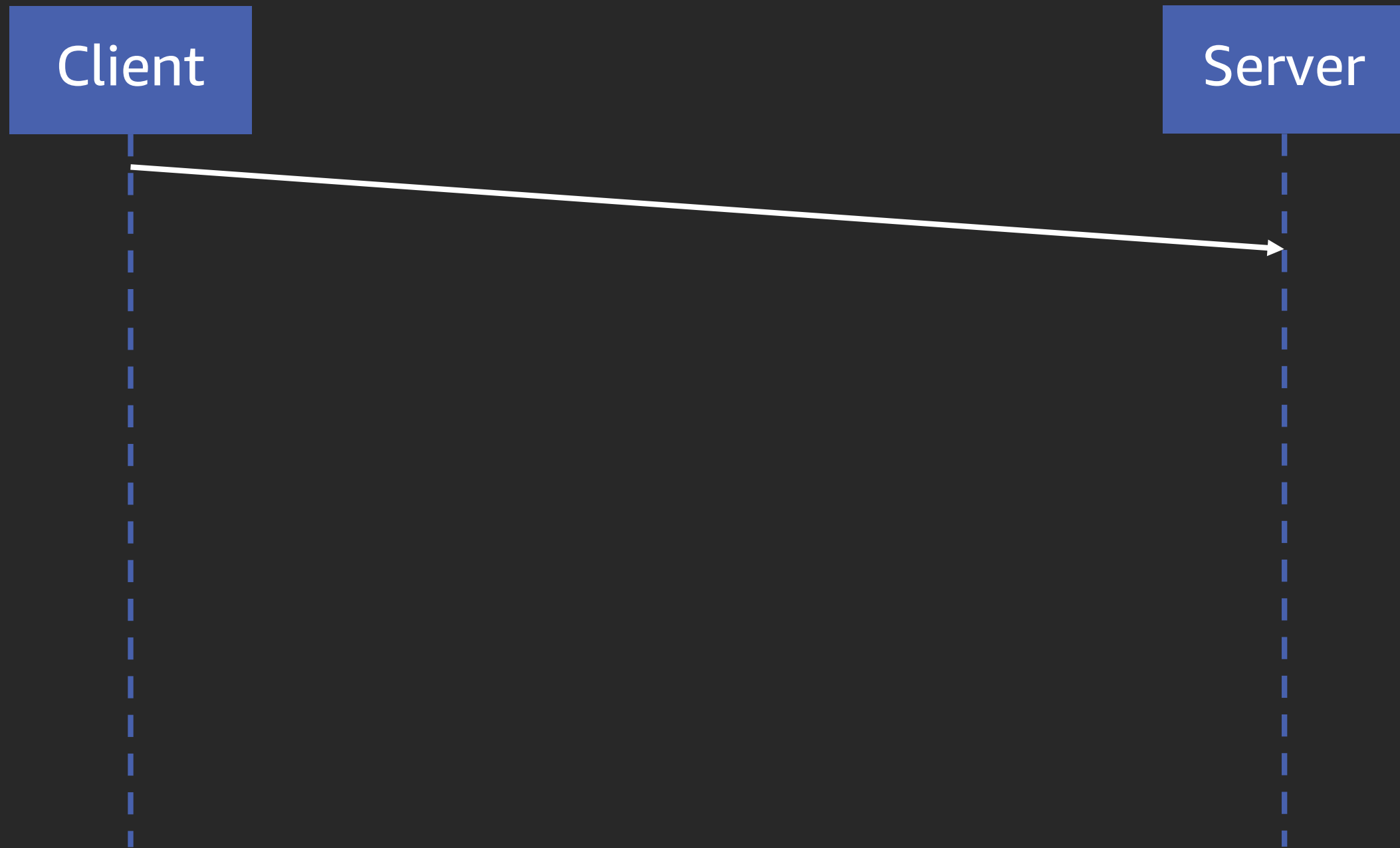Legend: Latency, Client Timeout

Goodput vs throughput,
with and without load shedding

# Don't waste work
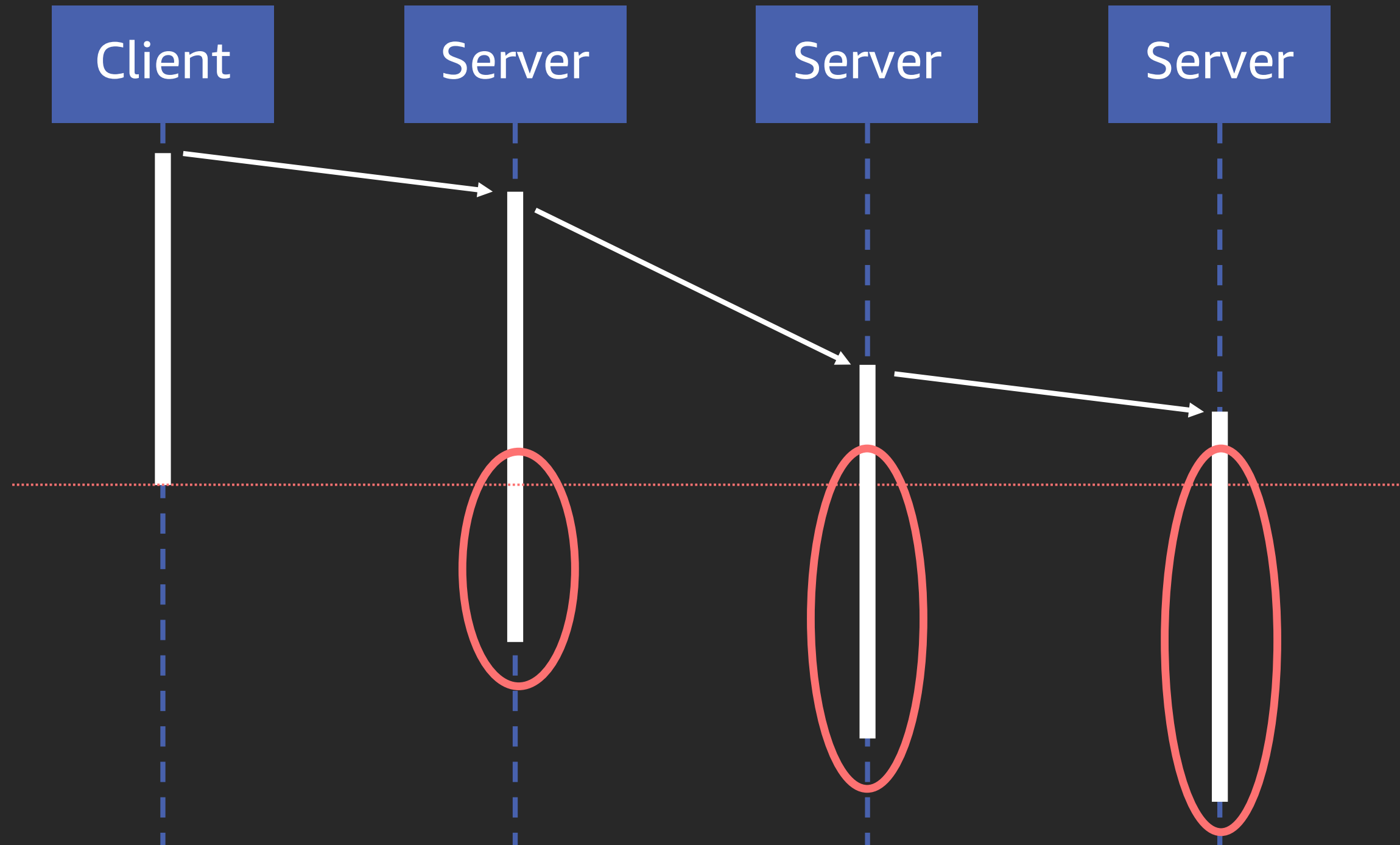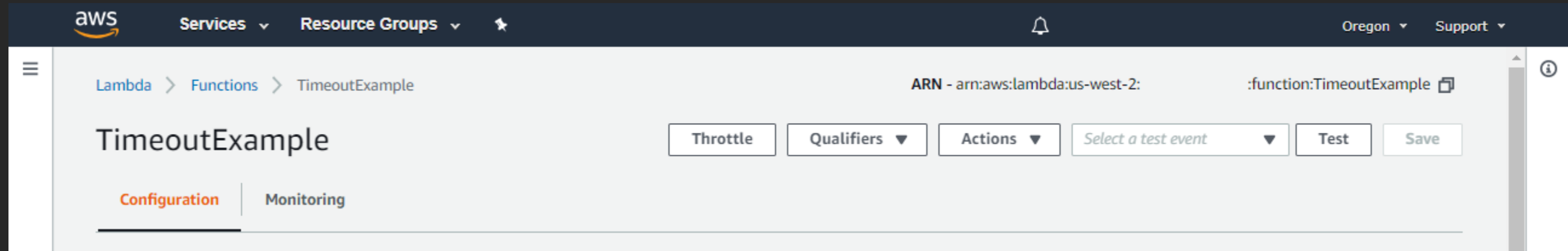
# Timeout wasting work

# Timeout wasting work

Client

Server

(client timeout)

# Timeout wasting work



(server keeps working)

# Layers of waste

# Server(less?)-side timeouts in Lambda

# Server-side timeouts



Client　　　　　　　　　　　　　　　　　　Server

(server stops after timeout)

(minimized waste)

# Better, still not great

## Goodput vs throughput



Y-axis: Goodput (TPS) — 0, 50, 100, 150, 200, 250

X-axis: Throughput (TPS) — 0, 50, 100, 150, 200, 250, 300, 350

Legend: Before · With server-side timeouts

# (Un?)Intended consequences?

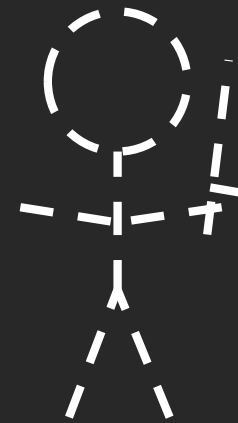Expensive requests?                    Slow dependency?
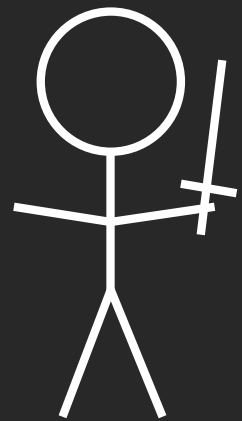
# Bounded work
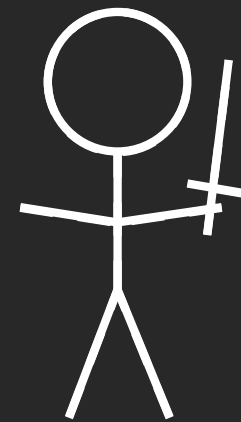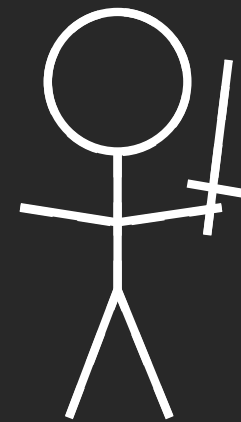
Input size validation

Pagination

Checkpointing

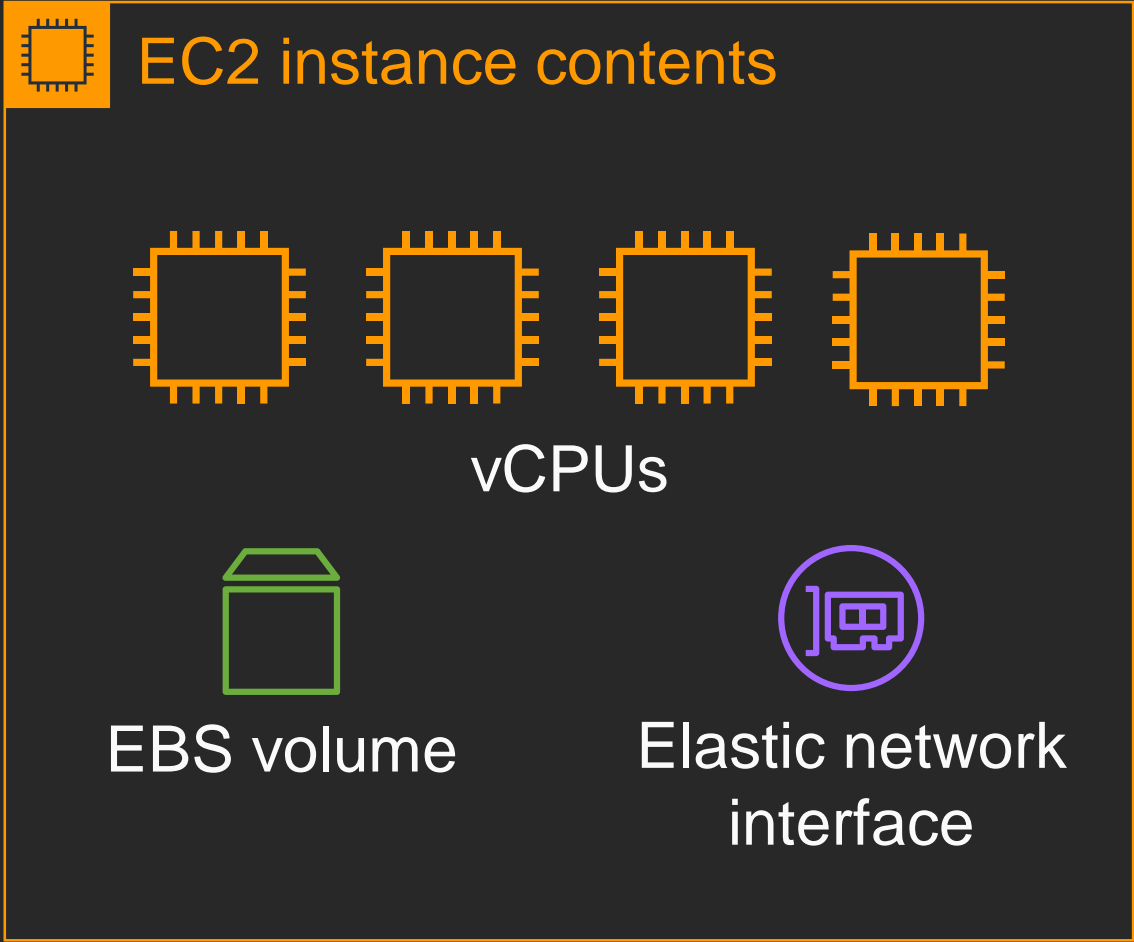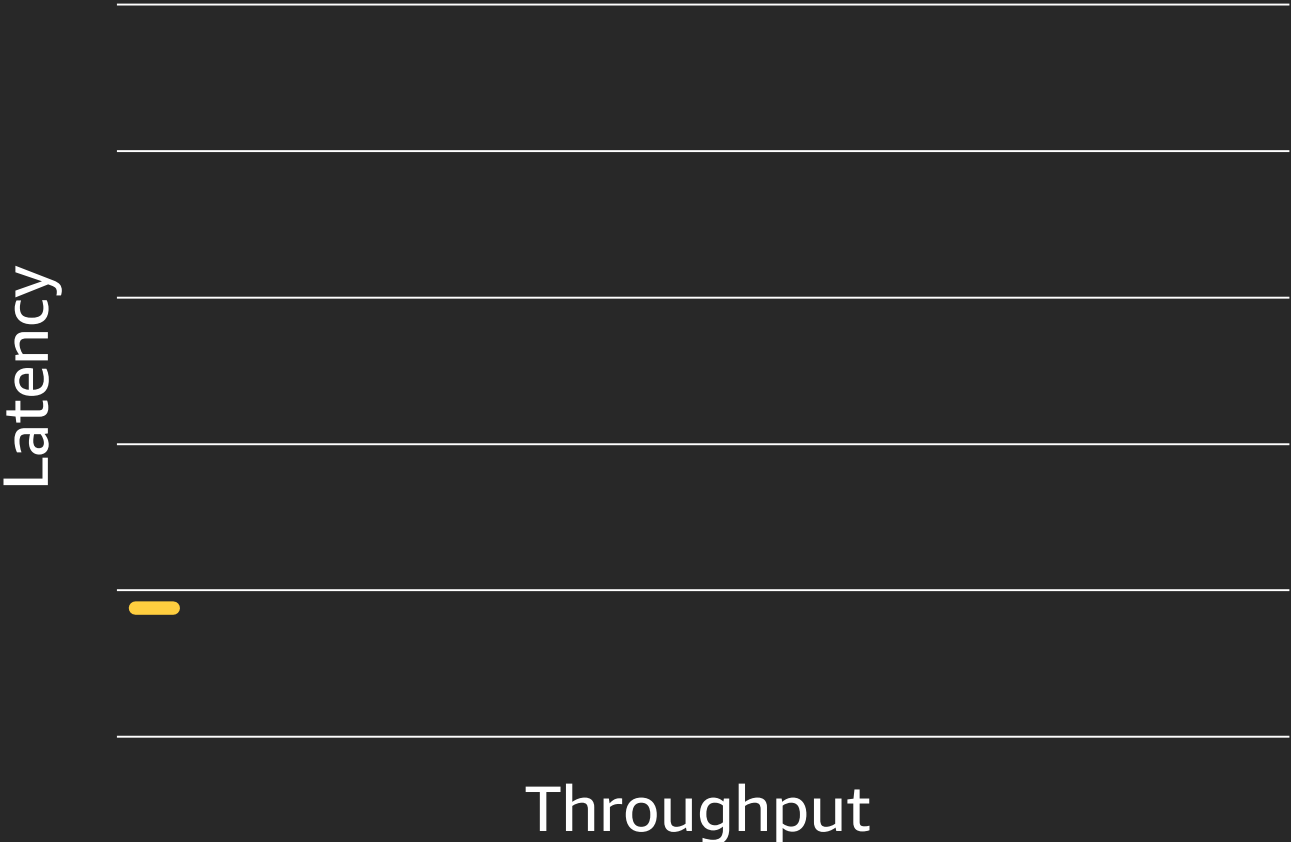# Checkpointing

# Checkpointing

# Checkpointing

# Bounded work

# Don't take on too much work

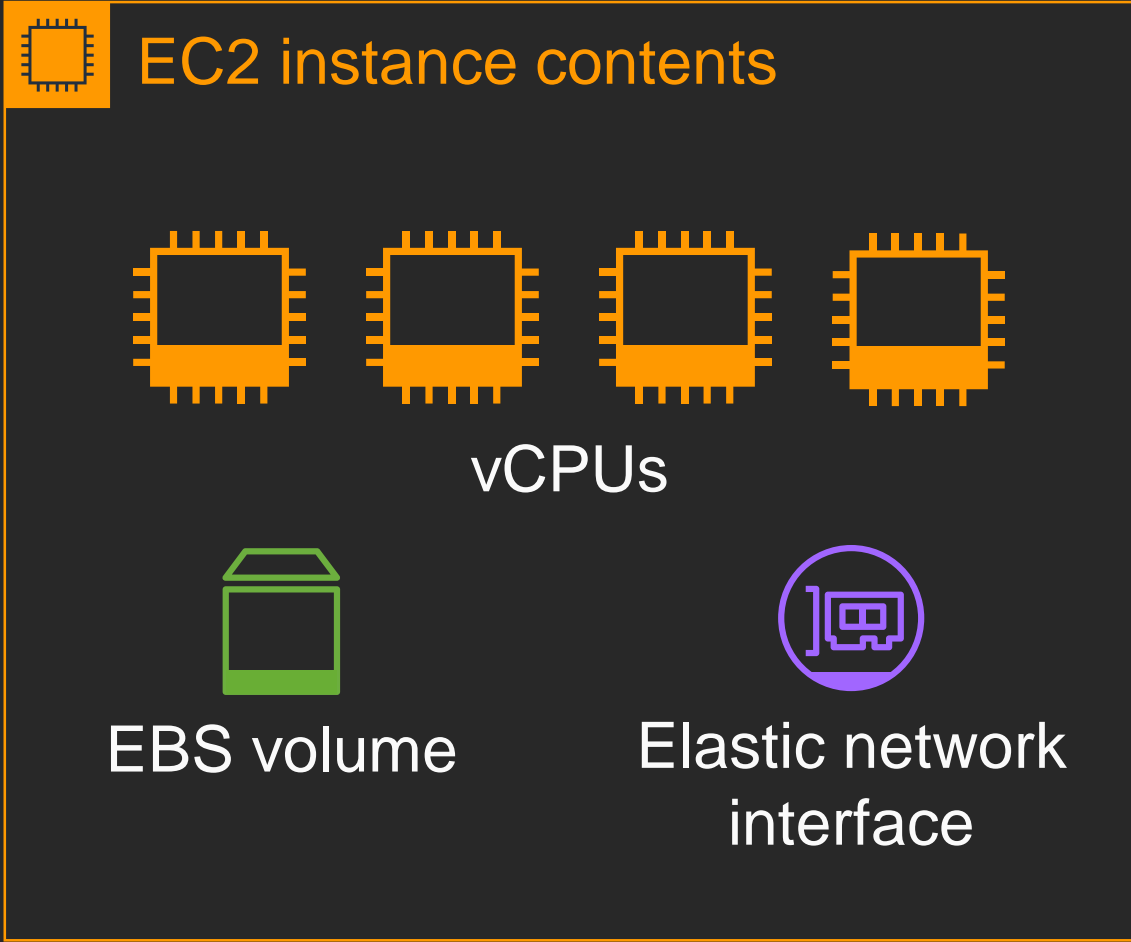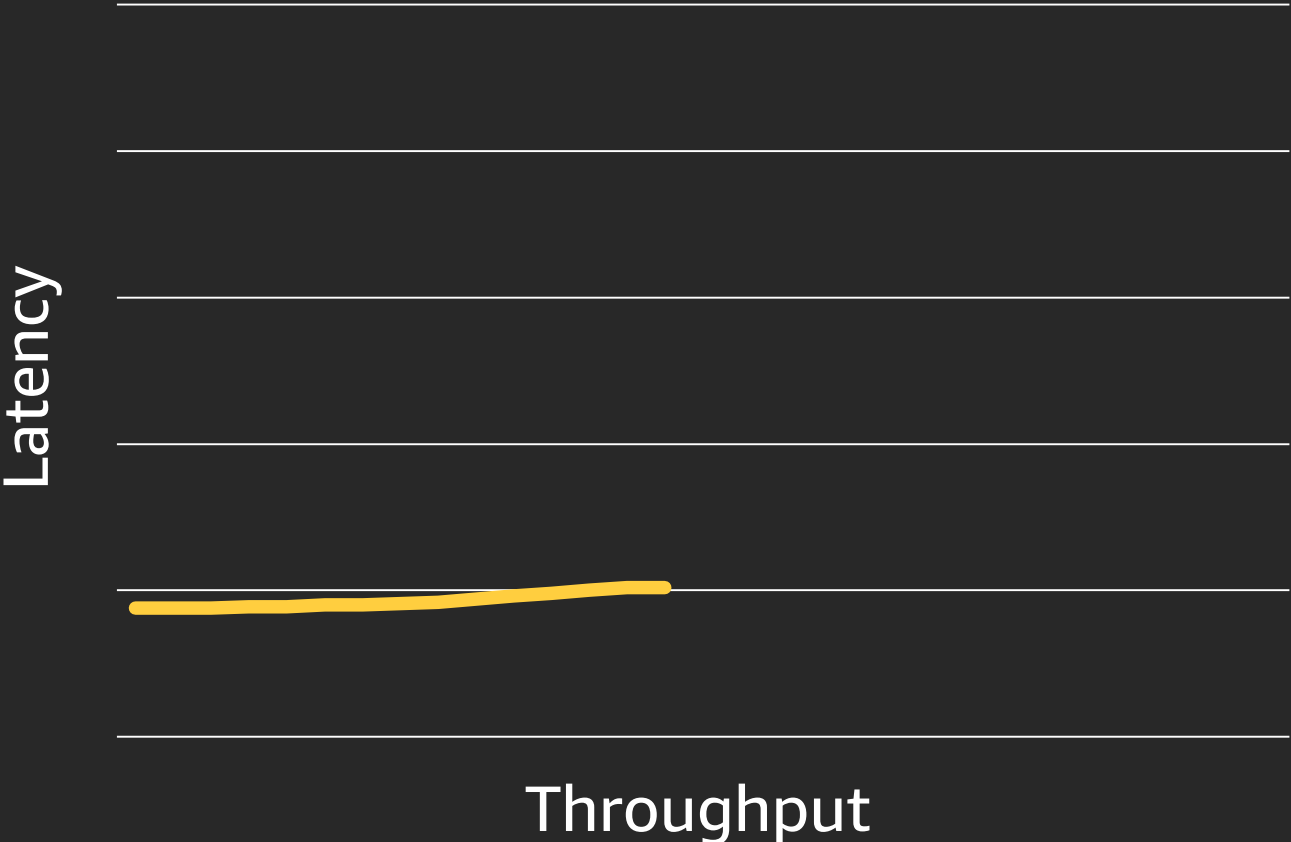# Degrading performance with load

## Latency vs throughput

Latency

Throughput

## EC2 instance contents

vCPUs

EBS volume

Elastic network interface

# Degrading performance with load

## Latency vs throughput

Latency

Throughput

## EC2 instance contents

vCPUs

EBS volume

Elastic network interface

# Degrading performance with load

## Latency vs throughput



(Chart: y-axis labeled "Latency", x-axis labeled "Throughput", showing a yellow curve that stays low and flat before rising sharply at high throughput.)

## EC2 instance contents

vCPUs

EBS volume

Elastic network interface

# Degrading performance with load

## Latency vs throughput

Latency

Throughput

## EC2 instance contents

vCPUs

EBS volume

Elastic network interface

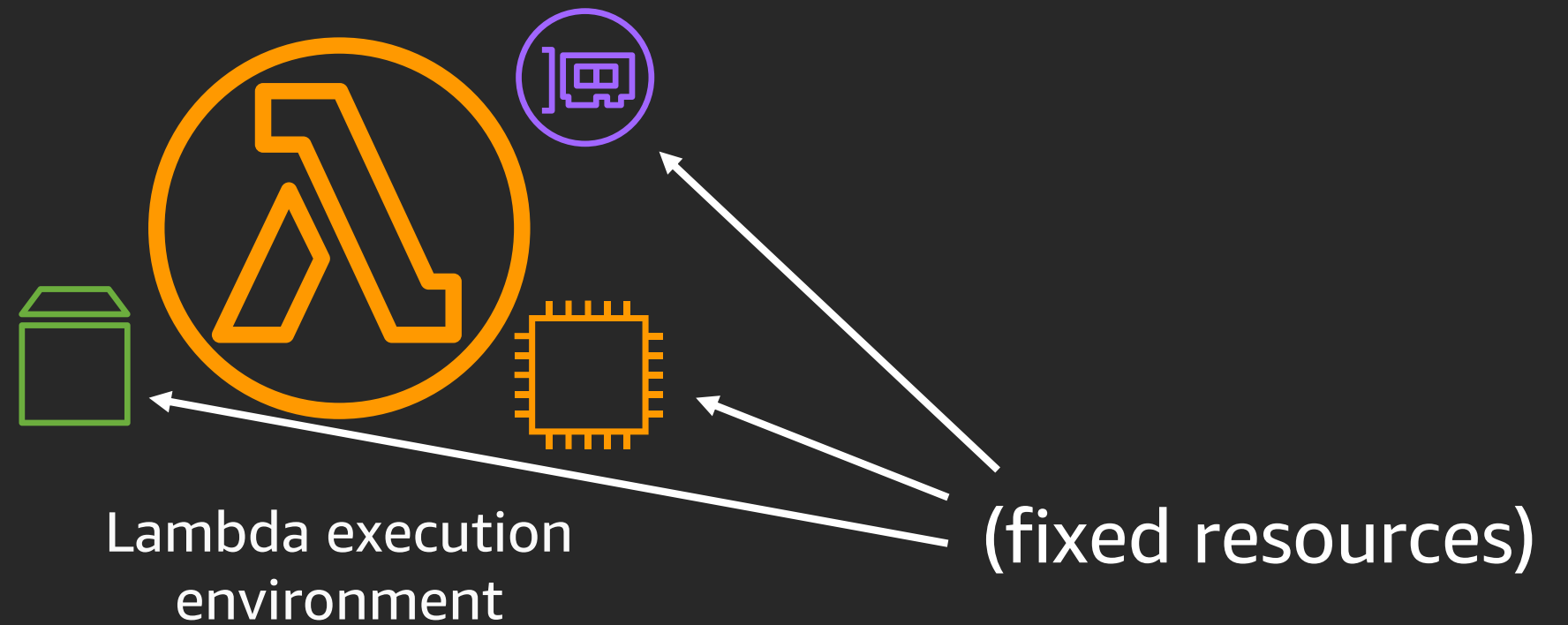# Fixed resources per unit of work



Lambda execution
environment

# Fixed resources per unit of work



Lambda execution
environment

(fixed resources)

# Fixed resources per unit of work



Request1    Request2

Lambda execution environment

Lambda execution environment

**(only work on one thing at a time)**

# Fixed resources per unit of work

Request1  Request2  Request3

Lambda execution environment

Lambda execution environment

Lambda execution environment

**(only work on one thing at a time)**

# Fixed resources per unit of work

Request4

Request2  Request3

Lambda execution
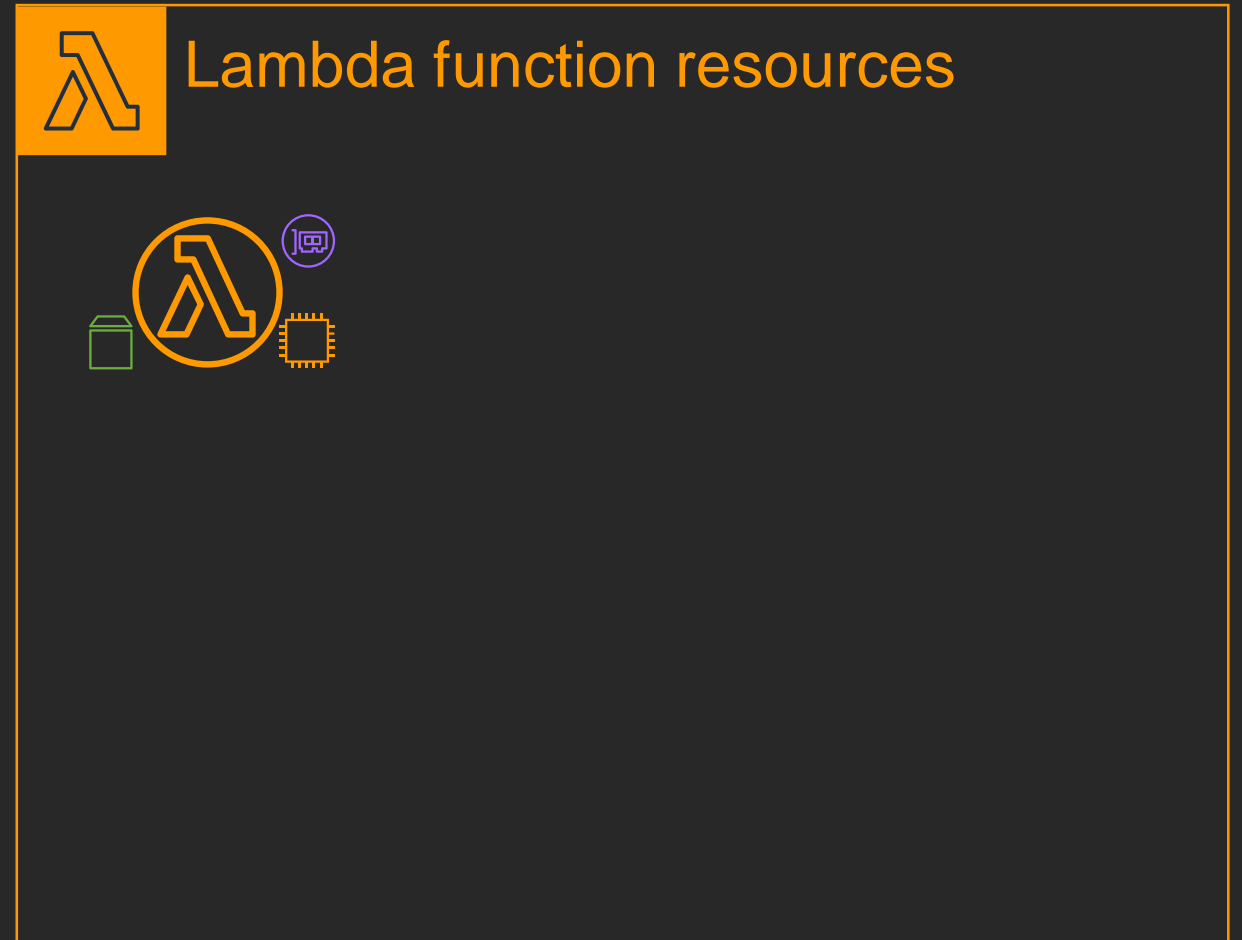environment

Lambda execution
environment

Lambda execution
environment

**(only work on one thing at a time)**

# Workload isolation means predictable performance

## Latency vs throughput



Latency

Throughput
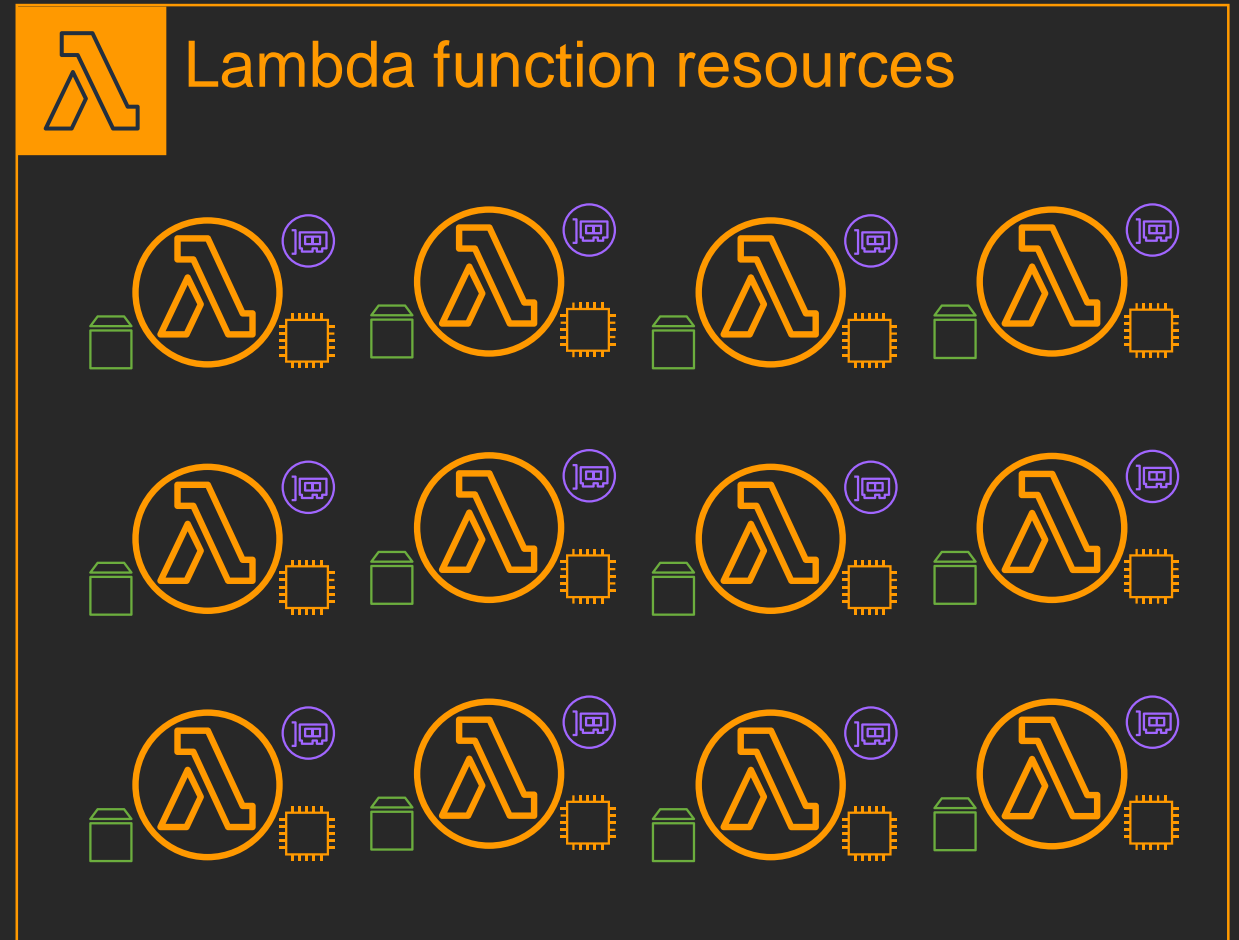
Lambda function resources

# Workload isolation means predictable performance
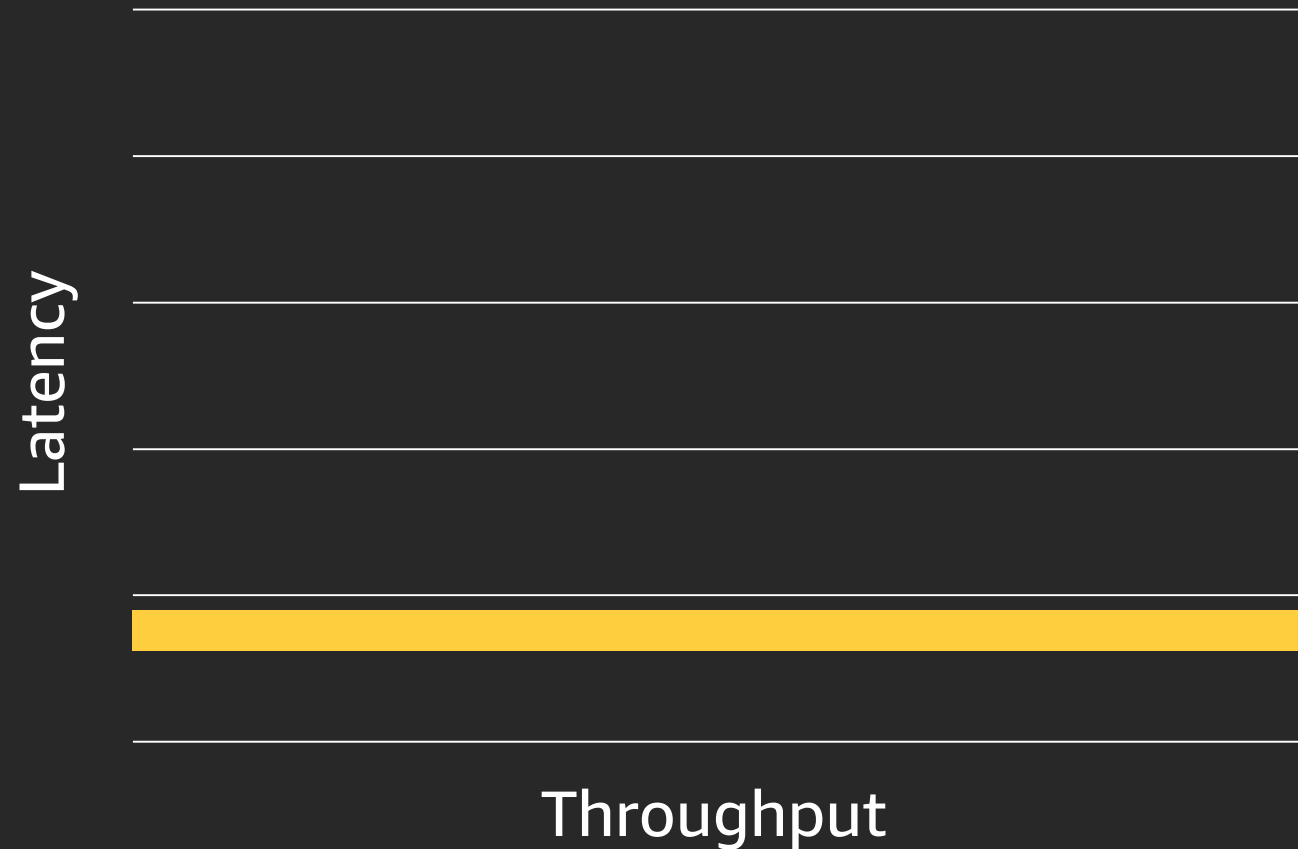
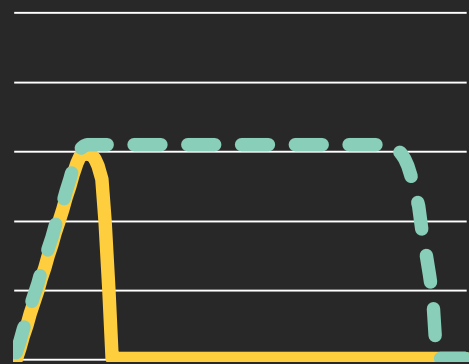## Latency vs throughput

Latency

Throughput



Lambda function resources

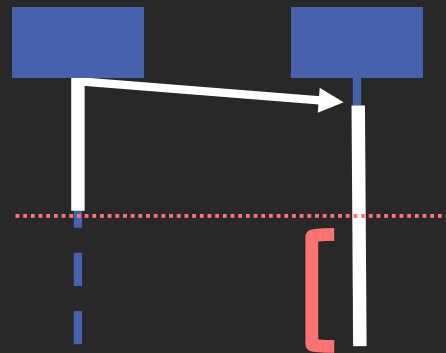# Workload isolation means predictable performance

# Don't take on too much work
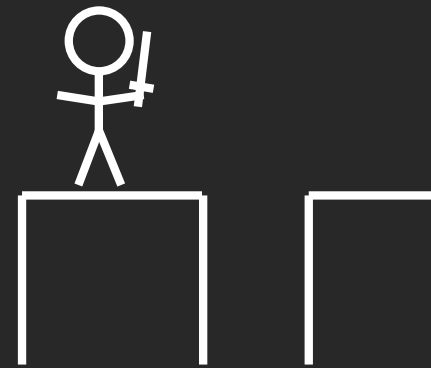
# Load shedding recap

**Reject excess load**

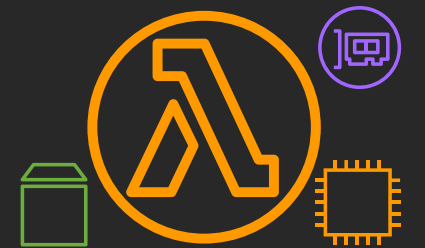Keep latency low for the work you take on

**Don't waste work**

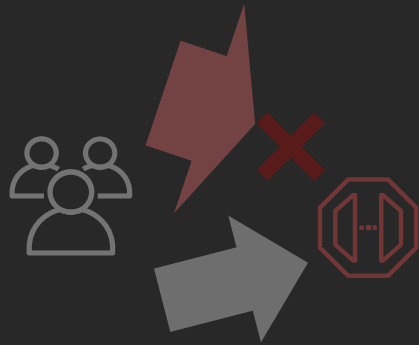Wasting work reduces throughput

**Do bounded work**

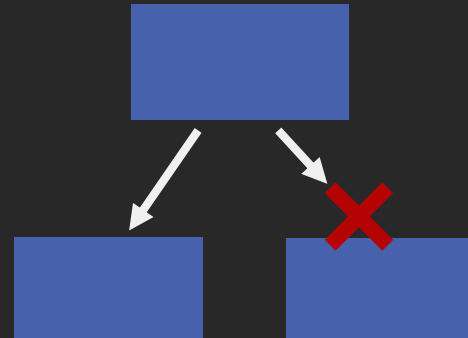Make scaling more predictable and testable

**Don't take on extra work**

Give the same resources to each request

# Chapter two

**Load shedding**

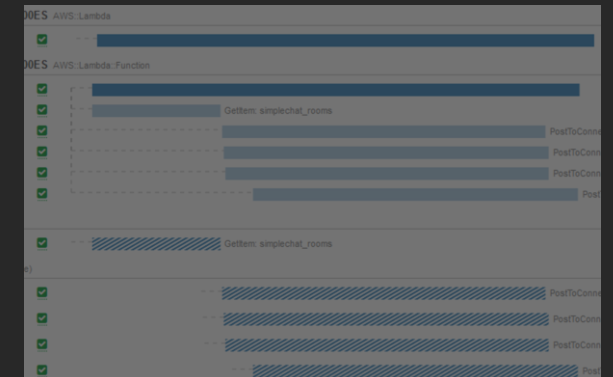Avoid brownout by rejecting excess load

**Dependency isolation**

**Prevent one dependency from affecting unrelated functionality**

**Avoiding queue backlogs**

Prevent a backlog from extending recovery time
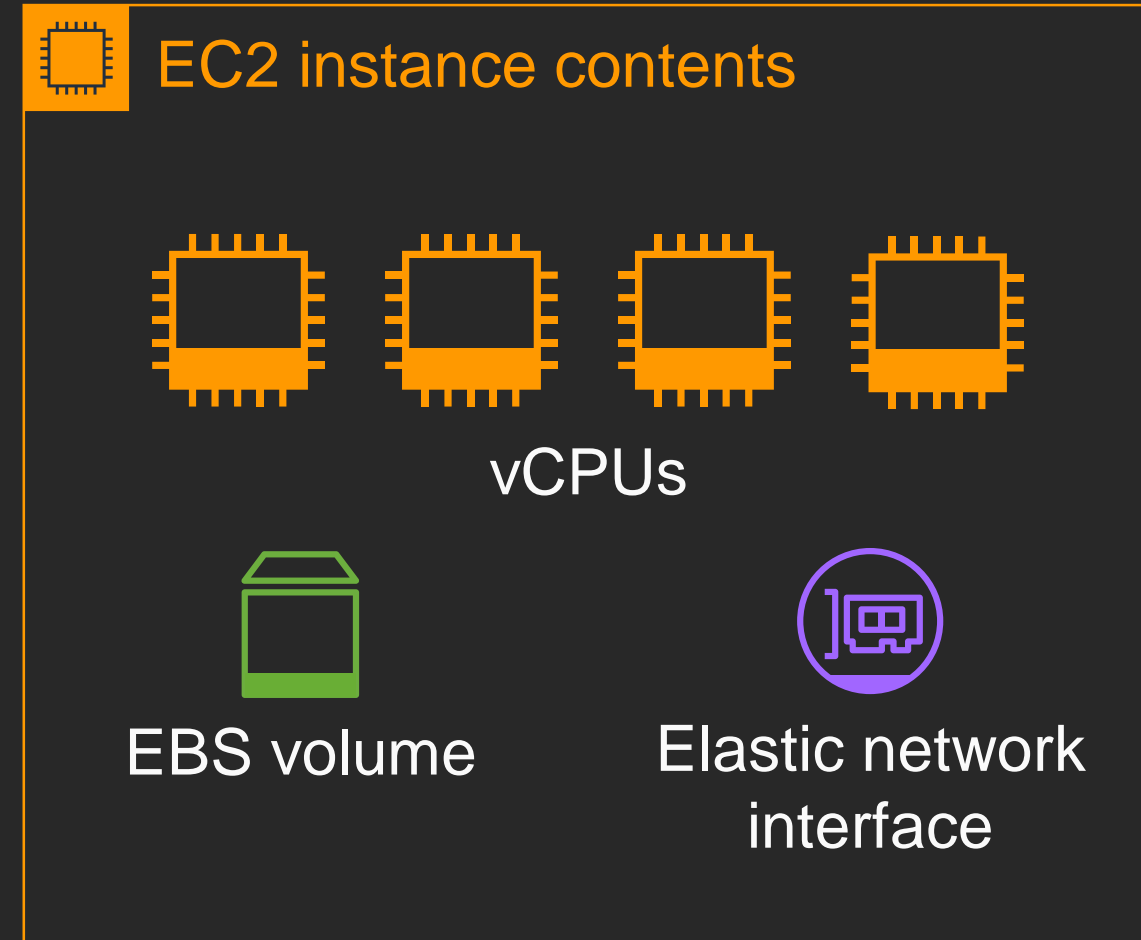
**Operating**

Quickly diagnose and mitigate issues

# Running out of concurrency

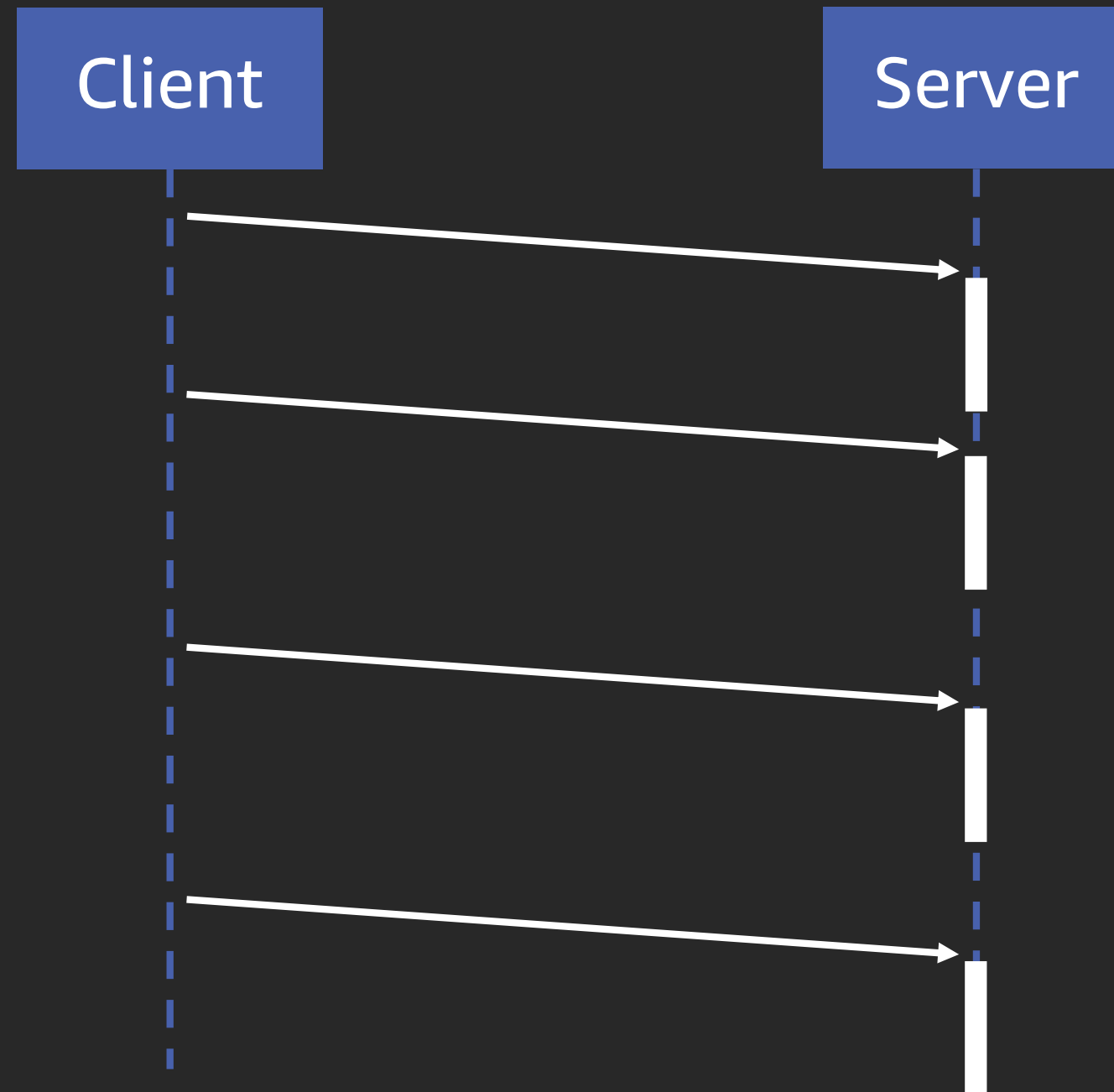# Application bottlenecks

Not pictured:

- Application thread pools
- File descriptors
- Ephemeral ports



EC2 instance contents

vCPUs

EBS volume

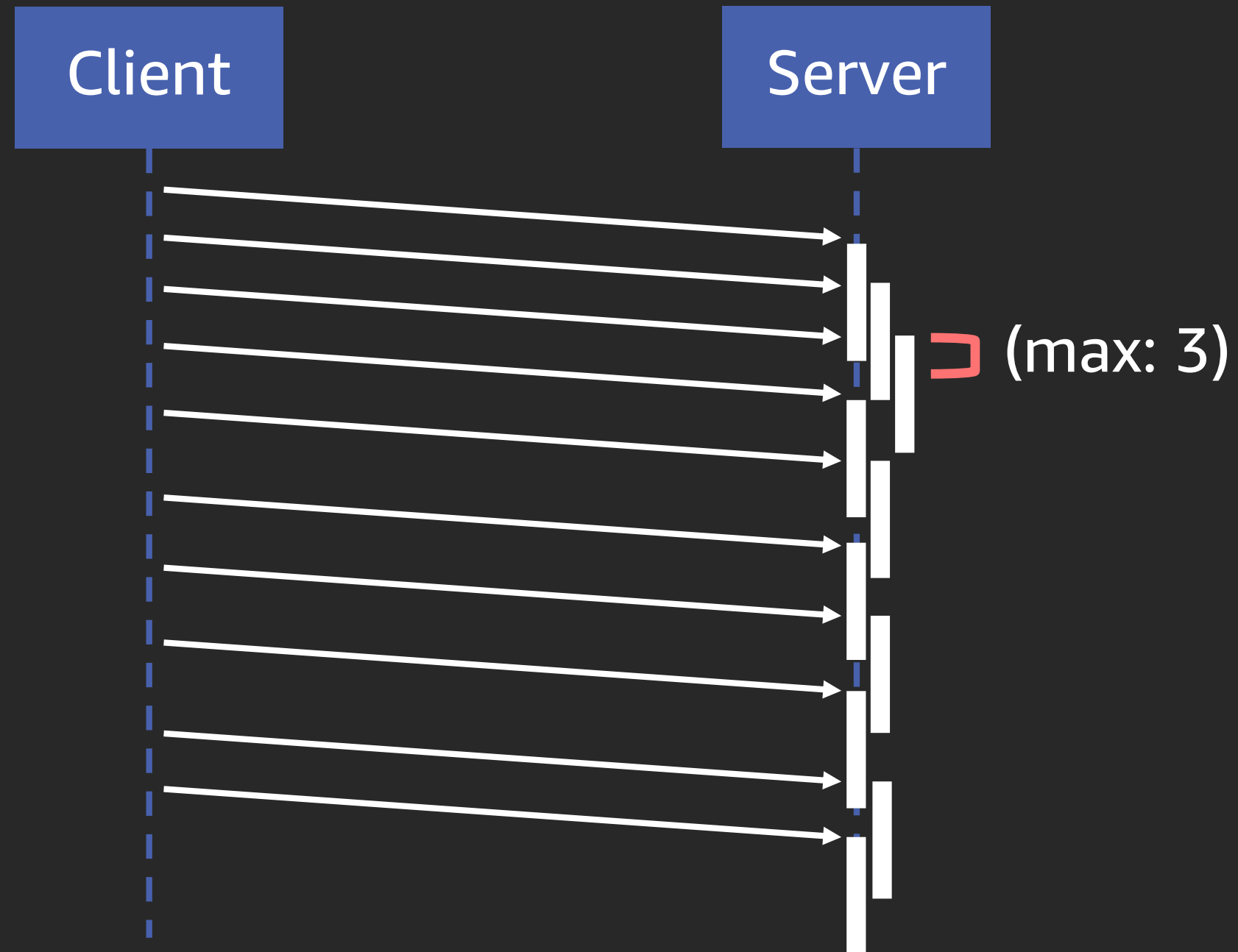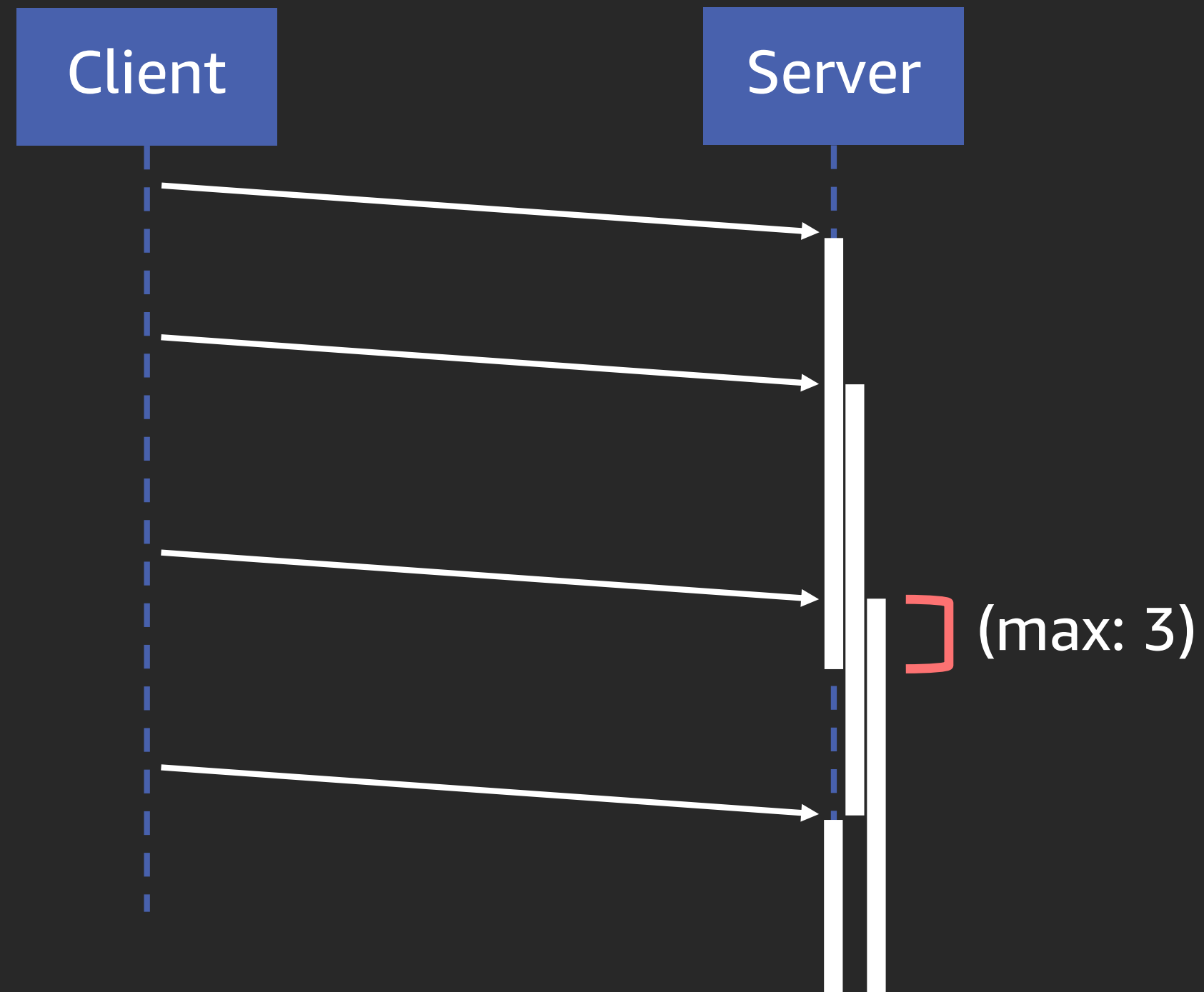Elastic network interface

# Concurrency, visualized

Concurrency with low request rate, low latency

# Concurrency increases with arrival rate

# Concurrency increases with latency



(max: 3)

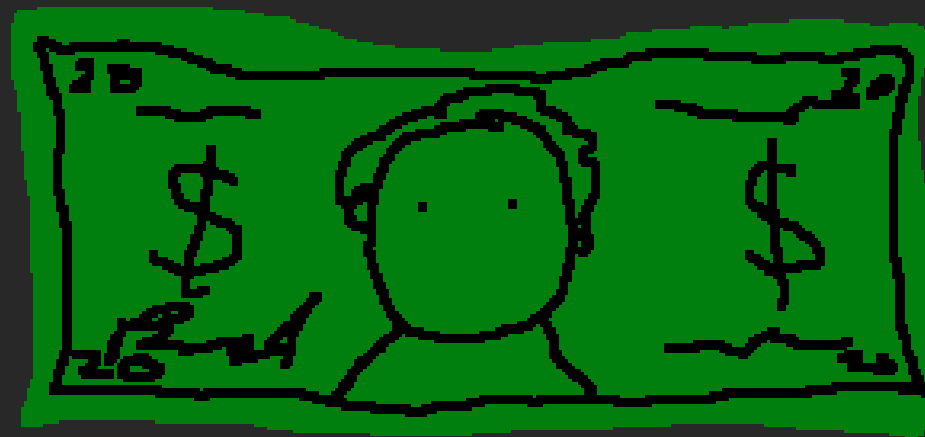# Little's Law

$$concurrency = arrival\ rate \ \times \ latency$$

$$concurrency = arrival\ rate \times latency$$

Cash

Check

# Dependency isolation within AWS Lambda

# Lambda architecture*

*currently

Lambda architecture: Cold invoke

# Lambda architecture: Warm invoke

# Warm invoke

# Cold invoke

# Concurrency from warm invokes



Max concurrency == 1

Flipping modes to cold invokes

Frontend   Worker manager   Placement

Max concurrency == 4

# Concurrency limits

**Frontend**

**Worker manager**

**Placement**

Concurrency limit == 4
(thread pool size,
file descriptors, etc)

# Approaching concurrency limits

# Exceeding concurrency limits

# Worker manager behavior

- Automatically scale
- Spread out increases in traffic to other worker managers
- Favor reusing existing sandboxes over creating new ones

# Worker manager behavior

- Automatically scale

- Spread out increases in traffic to other worker managers

- Favor reusing existing sandboxes over creating new ones

*warm start > cold start*

# Isolating concurrency

| Frontend | | Worker manager | | Placement |

Warm invokes only →

← Cold invokes only

# Isolating concurrency

# Isolating concurrency



Cold start retried on a different worker manager

# Isolating concurrency



Warm start availability unaffected

# Dependency isolation: Why?



**Isolate
unrelated APIs**

**Protect against
modal behavior**

# API isolation: How?



**Amazon API Gateway
per-API throttling**



**Lambda per-function
concurrency controls**

# Modal behavior protection: How?



CACHE

API Gateway throttling

CACHE

Per-function concurrency

# Chapter three

**Load shedding**

Avoid brownout by rejecting excess load

**Dependency isolation**

Prevent one dependency from affecting unrelated functionality

**Avoiding queue backlogs**

**Prevent a backlog from extending recovery time**

**Operating**

Quickly diagnose and mitigate issues

# Basic chat architecture



https://github.com/aws-samples/simple-websockets-chat-app

# Business as usual



Users → Endpoint →

**Chat service**

SQS queue → Lambda function → Chat tables

# Surge subsides, backlog remains



Users → Chat service [ SQS queue → Lambda function → Chat tables ]

# Ideal customer experience?

Queue backlogs are bad

# FIFO

## (**F**irst **I**n, **F**irst **O**ut)

Producer

Newest

(huge backlog slows *all* messages)

Oldest

Consumer

# LIFO

**(L**ast **I**n, **F**irst **O**ut**)**

Producer

Newest

Consumer

(huge backlog slows *old* messages)

Oldest

# Making FIFO behave LIFO

# Priority queues

Chat service

High-priority queue

(check here first)

Low-priority queue

(check here when "high" is empty)

# Priority is not known upfront

Chat service

(Enqueue)

High priority

Low priority

# Priority is not known upfront

Chat service

(Enqueue)

High priority

Low priority

# Move old messages to low-priority queue



Chat service

(Process)

(Receive)

(Re-enqueue) (Receive)

High priority

Low priority

# Time to Live (TTL)

IoT service

High-priority queue

(Just delete old messages)

# Backpressure (throttling)

Users → API Gateway → **Chat service**: Queue → Lambda function → Chat tables

# Backpressure (throttling)



Users → API Gateway → Queue → Lambda function → Chat tables

Chat service

# Backpressure (throttling)

**(throttled)**

Users → API Gateway → Queue → Chat service: Lambda function → Chat tables

**Chat service**

Priority queues + throttling: Best of both worlds?

# Shuffle-sharding

# Under the hood: Lambda async

# Lambda async

# Lambda async



Async queue

Frontend

Poller

# Lambda async

# Lambda async (this is not what happens)

# Queue per workload

Frontend

(many more queues)

Poller

# Queue per workload



Frontend

(many more queues)

Poller

# Polling is not free

# Shuffle-sharding

# Shuffle-sharding

(fixed number of N queues)

# Shuffle-sharding



(fixed number of N queues)

(map each workload to K queues)

# Shuffle-sharding: Enqueue



(find emptier)

Frontend

Poller

# Shuffle-sharding: Enqueue

# Shuffle-sharding: Busy workload

(hot workload)

Frontend

Poller

# Shuffle-sharding: Busy workload



(find emptier)

Frontend

Poller

# Shuffle-sharding: Busy workload



(enqueue)

Frontend

Poller

# Shuffle-sharding: Magical resource isolation



$$\star \approx \frac{shardsize!}{|nodes|!}$$

# Shuffle-sharding: Magical resource isolation

**Nodes** = 8
**Shard size** = 2

| Overlap | % workloads |
|---------|-------------|
| 0 | 53.6% |
| 1 | 42.8% |
| 2 | 3.6% |

# Shuffle-sharding: Magical resource isolation

**Nodes = 100**
**Shard size = 5**

| Overlap | % workloads |
|---------|-------------|
| 0 | 77% |
| 1 | 21% |
| 2 | 1.8% |
| 3 | 0.06% |
| 4 | 0.0006% |
| 5 | 0.0000013% |

# Isolate busy workloads

(isolate busy workload)

# Avoiding queue backlogs

- Backlogs build quickly, introduce modes

# Avoiding queue backlogs

- Backlogs build quickly, introduce modes
- Auto Scaling and Lambda react quickly

Amazon EC2
Auto Scaling

AWS
Lambda

# Avoiding queue backlogs

- Backlogs build quickly, introduce modes

- Auto Scaling and Lambda react quickly

- Priority queueing emulates LIFO

Chat service

High priority    Low priority

# Avoiding queue backlogs

- Backlogs build quickly, introduce modes
- Auto Scaling and Lambda react quickly
- Priority queueing emulates LIFO
- Move old messages to low priority

Chat service

High priority    Low priority

# Avoiding queue backlogs

- Backlogs build quickly, introduce modes

- Auto Scaling and Lambda react quickly

- Priority queueing emulates LIFO

- Move old messages to low priority

- Message TTLs for stale information

Chat service

High priority

(delete)

# Avoiding queue backlogs

- Backlogs build quickly, introduce modes
- Auto Scaling and Lambda react quickly
- Priority queueing emulates LIFO
- Move old messages to low priority
- Message TTLs for stale information
- Apply backpressure

(throttled)

Queue

# Avoiding queue backlogs

- Backlogs build quickly, introduce modes
- Auto Scaling and Lambda react quickly
- Priority queueing emulates LIFO
- Move old messages to low priority
- Message TTLs for stale information
- Apply backpressure
- Surge queue excess traffic

Surge queue

Warm queue

# Avoiding queue backlogs

- Backlogs build quickly, introduce modes
- Auto Scaling and Lambda react quickly
- Priority queueing emulates LIFO
- Move old messages to low priority
- Message TTLs for stale information
- Apply backpressure
- Surge queue excess traffic
- Shuffle-sharding for isolation

# Chapter four

Load shedding

Avoid brownout by
rejecting excess load

Dependency
isolation

Prevent one
dependency from
affecting unrelated
functionality

Avoiding queue
backlogs

Prevent a backlog
from extending
recovery time

**Operating**

**Quickly diagnose
and mitigate issues**

# AWS X-Ray tracing

## Traces › Details

**Timeline** | **Raw data**

| Method | Response | Duration | Age | | ID |
|--------|----------|----------|-----|---|-----|
| -- | 200 | 309 ms | 5.3 min (2019-03-16 22:59:06 UTC) | | 1-5c8d7fba-b3b2ae429da8149f3e47df66 |

| Name | Res. | Duration | Status | 0.0ms | 50ms | 100ms | 150ms | 200ms | 250ms | 300ms | 350ms |
|------|------|----------|--------|-------|------|-------|-------|-------|-------|-------|-------|

▼ **serverlessrepo-simple-websocke-SendMessageFunction-R9QXLHH800ES** AWS::Lambda

| serverlessrepo-simple-websocke-SendMessageFun | 200 | 306 ms | ☑ | | | | | | | | |

▼ **serverlessrepo-simple-websocke-SendMessageFunction-R9QXLHH800ES** AWS::Lambda::Function

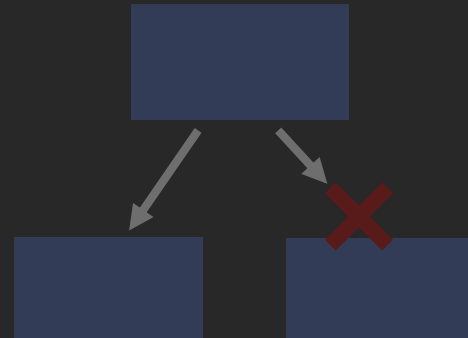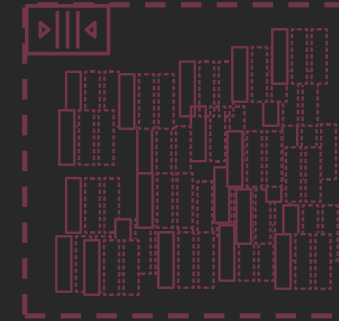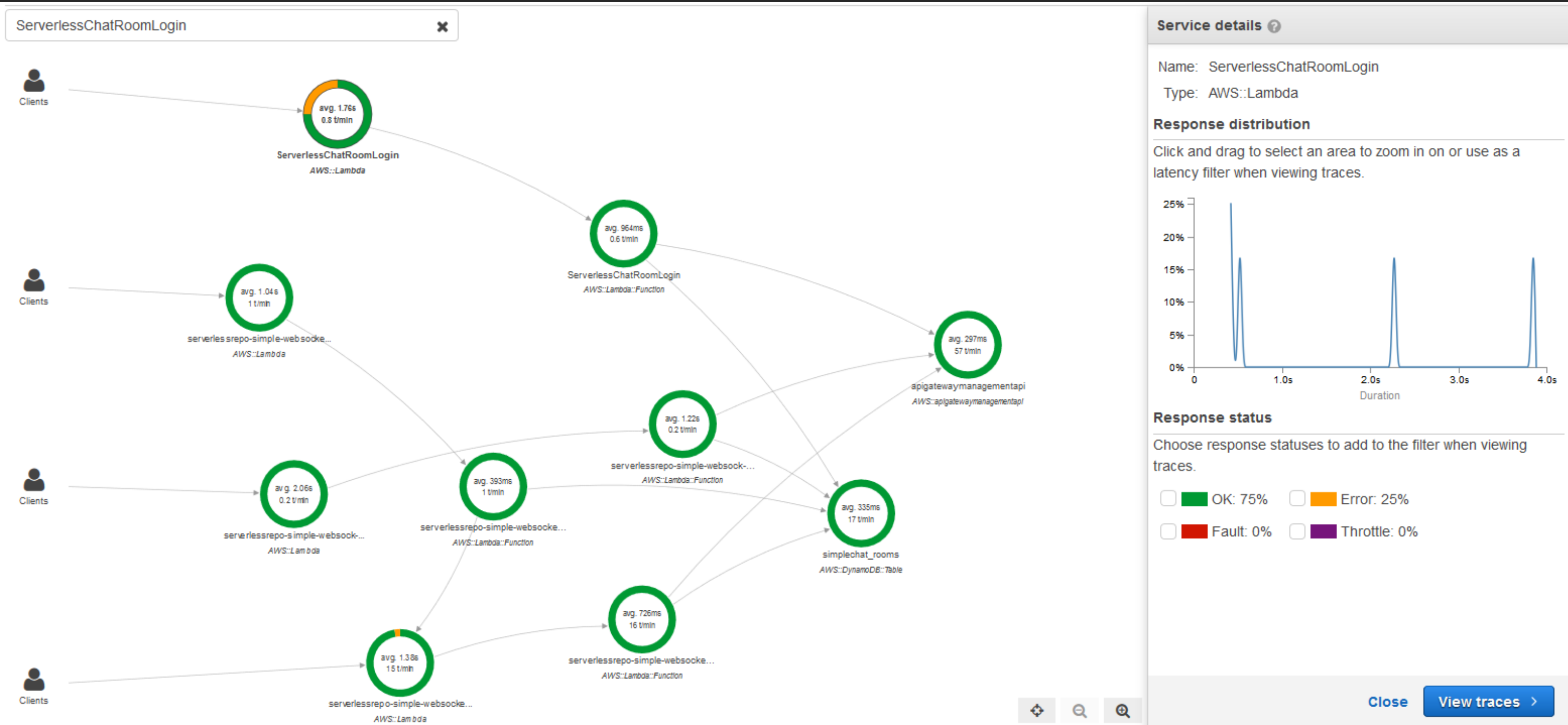| serverlessrepo-simple-websocke-SendMessageFun | - | 300 ms | ☑ | | | | | | | | |
| DynamoDB | 200 | 79.0 ms | ☑ | | GetItem: simplechat_rooms | | | | | | |
| apigatewaymanagementapi | 200 | 199 ms | ☑ | | | | | PostToConnection | | | |
| apigatewaymanagementapi | 200 | 200 ms | ☑ | | | | | PostToConnection | | | |
| apigatewaymanagementapi | 200 | 199 ms | ☑ | | | | | PostToConnection | | | |
| apigatewaymanagementapi | 200 | 200 ms | ☑ | | | | | PostToConnection | | | |

▼ **DynamoDB** AWS::DynamoDB::Table (Client Response)

| serverlessrepo-simple-websocke-SendMessageFun | 200 | 79.0 ms | ☑ | | GetItem: simplechat_rooms | | | | | | |

▼ **apigatewaymanagementapi** AWS::apigatewaymanagementapi (Client Response)

| serverlessrepo-simple-websocke-SendMessageFun | 200 | 199 ms | ☑ | | | | | PostToConnection | | | |
| serverlessrepo-simple-websocke-SendMessageFun | 200 | 200 ms | ☑ | | | | | PostToConnection | | | |
| serverlessrepo-simple-websocke-SendMessageFun | 200 | 199 ms | ☑ | | | | | PostToConnection | | | |

# Amazon CloudWatch insights

# Amazon CloudWatch insights

```
fields @timestamp, @duration, date_floor(@timestamp, 1s)
| filter @duration > 0
| stats avg(@duration), count(*), (avg(@duration) / 1000.0) * count(*) as concurrency by bin(1s)
| sort concurrency desc
```

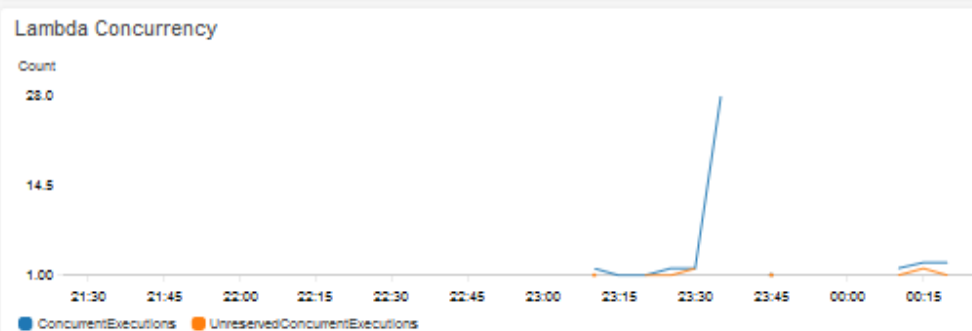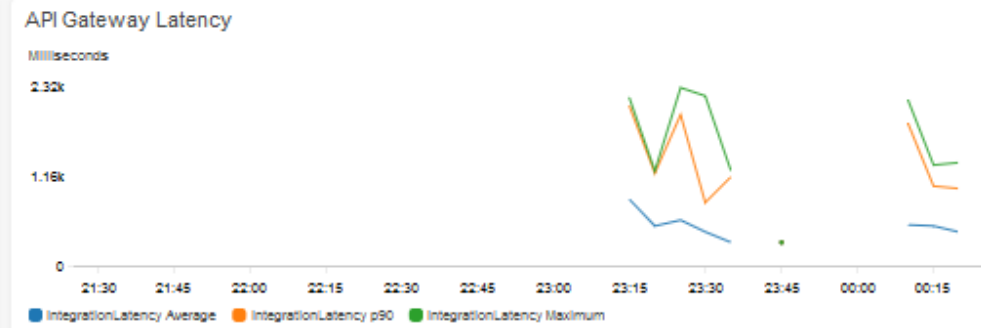| # | bin(1s) | avg(@duration) | count(*) | concurrency |
|---|---------|----------------|----------|-------------|
| 1 | 2019-03-14 17:26:17.000 | 3003.174 | 471 | 1414.4949 |
| 2 | 2019-03-14 17:26:18.000 | 3003.1106 | 351 | 1054.0918 |
| 3 | 2019-03-14 17:25:17.000 | 2400.892 | 65 | 156.058 |
| 4 | 2019-03-14 17:26:16.000 | 2818.4718 | 44 | 124.0128 |
| 5 | 2019-03-14 17:26:20.000 | 3002.2058 | 40 | 120.0882 |
| 6 | 2019-03-14 17:25:54.000 | 1878.0819 | 52 | 97.6603 |
| 7 | 2019-03-14 17:26:04.000 | 1903.5012 | 49 | 93.2716 |
| 8 | 2019-03-14 17:26:15.000 | 1868.6393 | 46 | 85.9574 |
| 9 | 2019-03-14 17:25:18.000 | 1315.5267 | 52 | 68.4074 |
| 10 | 2019-03-14 17:26:05.000 | 3003.1705 | 19 | 57.0602 |
| 11 | 2019-03-14 17:26:19.000 | 3002.1027 | 15 | 45.0315 |
| 12 | 2019-03-14 17:25:55.000 | 2231.6844 | 16 | 35.707 |
| 13 | 2019-03-14 17:25:53.000 | 1072.0553 | 32 | 34.3058 |
| 14 | 2019-03-14 17:25:16.000 | 1811.7122 | 18 | 32.6108 |
| 15 | 2019-03-14 17:26:14.000 | 1303.0379 | 24 | 31.2729 |

## API Gateway Volume

Count

## API Gateway Errors

Count

## API Gateway Latency

Milliseconds

## Lambda Concurrency

Count

## Lambda Volume

Count

## DynamoDB Capacity

Count

## Lambda SendMessage Volume

Count

## Lambda SendMessage Latency

Milliseconds

## Lambda SendMessage Concurrency

Count

## Lambda Login Volume

Count

## Lambda Login Latency

Milliseconds

Log group: /aws/lambda/serverlessrepo-simple-websocke-SendMessageFunction-R9QXLHH800ES

| # | @timestamp | @message |
|---|---|---|
| 1 | 2019-03-16 00:13:21.730 | REPORT RequestId: b3a09f85-5fd2-46f9-a430-7a22aa9e57d8 Duration: 166.70 ms Billed Duration: 200 ms Memory Size: 256 MB Max Memory Used: 90 MB |
| 2 | 2019-03-16 00:13:21.730 | END RequestId: b3a09f85-5fd2-46f9-a430-7a22aa9e57d8 |
| 3 | 2019-03-16 00:13:21.600 | 2019-03-16T00:13:21.600Z b3a09f85-5fd2-46f9-a430-7a22aa9e57d8 endpoint: 5ijpka7ch1.execute-api.us-west-2.amazonaws.com/dev |
| 4 | 2019-03-16 00:13:21.600 | 2019-03-16T00:13:21.600Z b3a09f85-5fd2-46f9-a430-7a22aa9e57d8 Version: 899, Aliases: Set { wrapperName: 'Set', values: [ '$', 'David', 'David1' ], type: 'String' }, Connections: { 'wn7BrdYsvMsCZvg+': { Alias: 'David' }, 'wn7i8eKsvMsAcuA+': { Alias: '$' }, 'wn8bdsQMPMcAcuA+': { Alias: 'David1' }, 'wn8i9F22vMsCtDQ+': {} }, ChatRoomId: 'default' } |
| 5 | 2019-03-16 00:13:21.600 | 2019-03-16T00:13:21.600Z b3a09f85-5fd2-46f9-a430-7a22aa9e57d8 Returned item |
| 6 | 2019-03-16 00:13:21.600 | 2019-03-16T00:13:21.600Z b3a09f85-5fd2-46f9-a430-7a22aa9e57d8 {"Type":"Message","Message":"anonymous","Alias":"$"} |
| 7 | 2019-03-16 00:13:21.563 | 2019-03-16T00:13:21.563Z b3a09f85-5fd2-46f9-a430-7a22aa9e57d8 Starting |
| 8 | 2019-03-16 00:13:21.562 | START RequestId: b3a09f85-5fd2-46f9-a430-7a22aa9e57d8 Version: $LATEST |
| 9 | 2019-03-16 00:13:18.790 | REPORT RequestId: 8se2f2fb-557c-449f-9ebd-c4656fbf42d8 Duration: 222.40 ms Billed Duration: 300 ms Memory Size: 256 MB Max Memory Used: 90 MB |
| 10 | 2019-03-16 00:13:18.790 | END RequestId: 8se2f2fb-557c-449f-9ebd-c4656fbf42d8 |

# Amazon CloudWatch Contributor Insights

## MessagesBySender

Show: Top 10 contributors | 1 minute period | Order by Sum

2019-11-29 (19:56:00) - 2019-11-29 (20:10:21)

Stacked area

6 unique contributors



1. David　2. Laura　3. Katie　4. Ben　5. Kyle　6. Sarah

| # | $.Alias | SampleCount | |
|---|---------|-------------|---|
| 1 | David | 154 | |
| 2 | Laura | 70 | |
| 3 | Katie | 47 | |
| 4 | Ben | 42 | |
| 5 | Kyle | 41 | |
| 6 | Sarah | 37 | |

# Amazon CloudWatch ServiceLens



Overall system health

Drill down to a specific request

# Recap: Resilient systems

aws

# Don't take on too much work



## Goodput vs throughput

# Reject excess work

Goodput vs throughput



Goodput (TPS) [y-axis]

Throughput (TPS) [x-axis]

— Goodput    - - With Load Shedding

# Workload isolation means predictable performance

## Latency vs throughput



Latency

Throughput

Request1

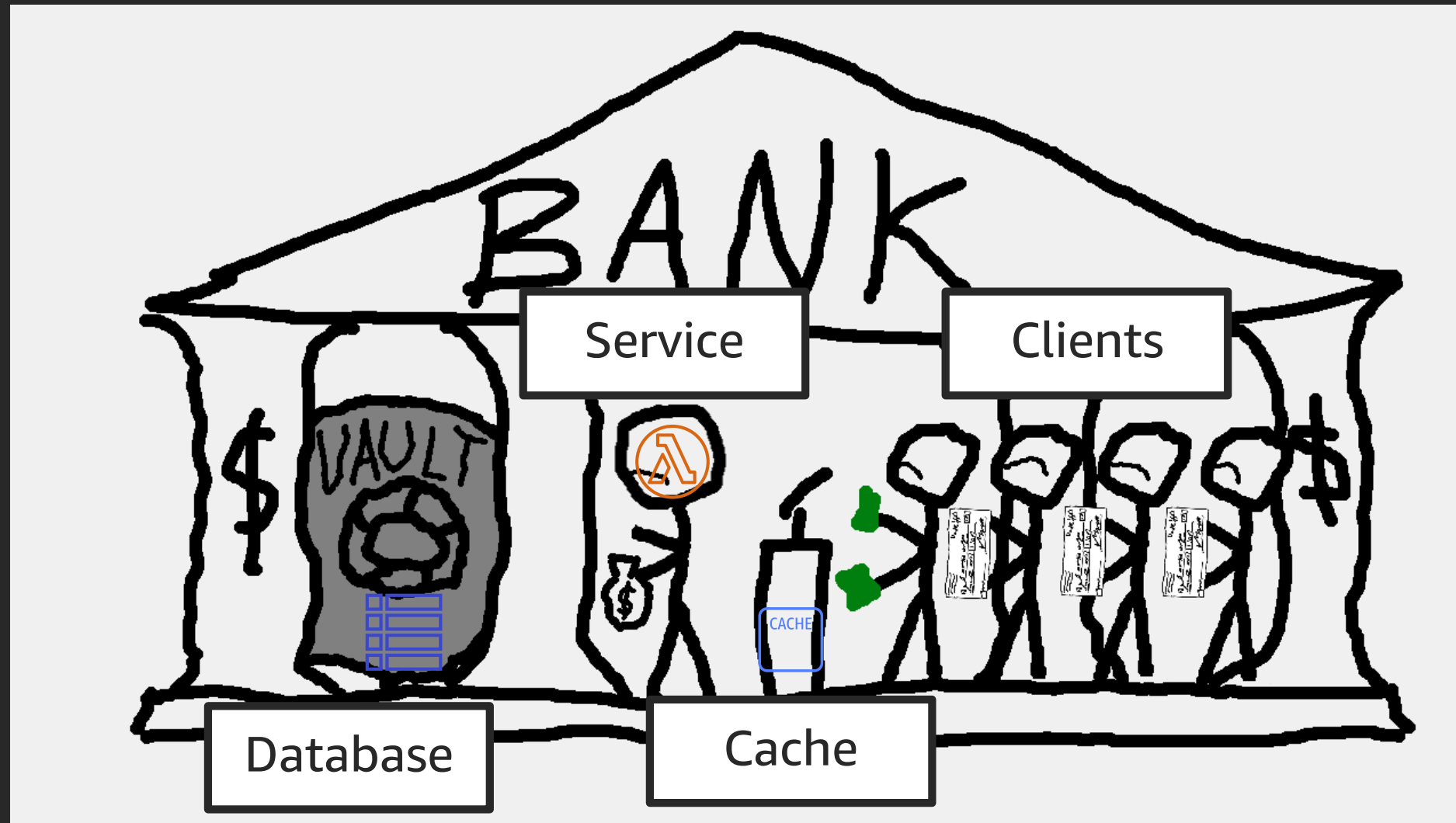Lambda execution environment

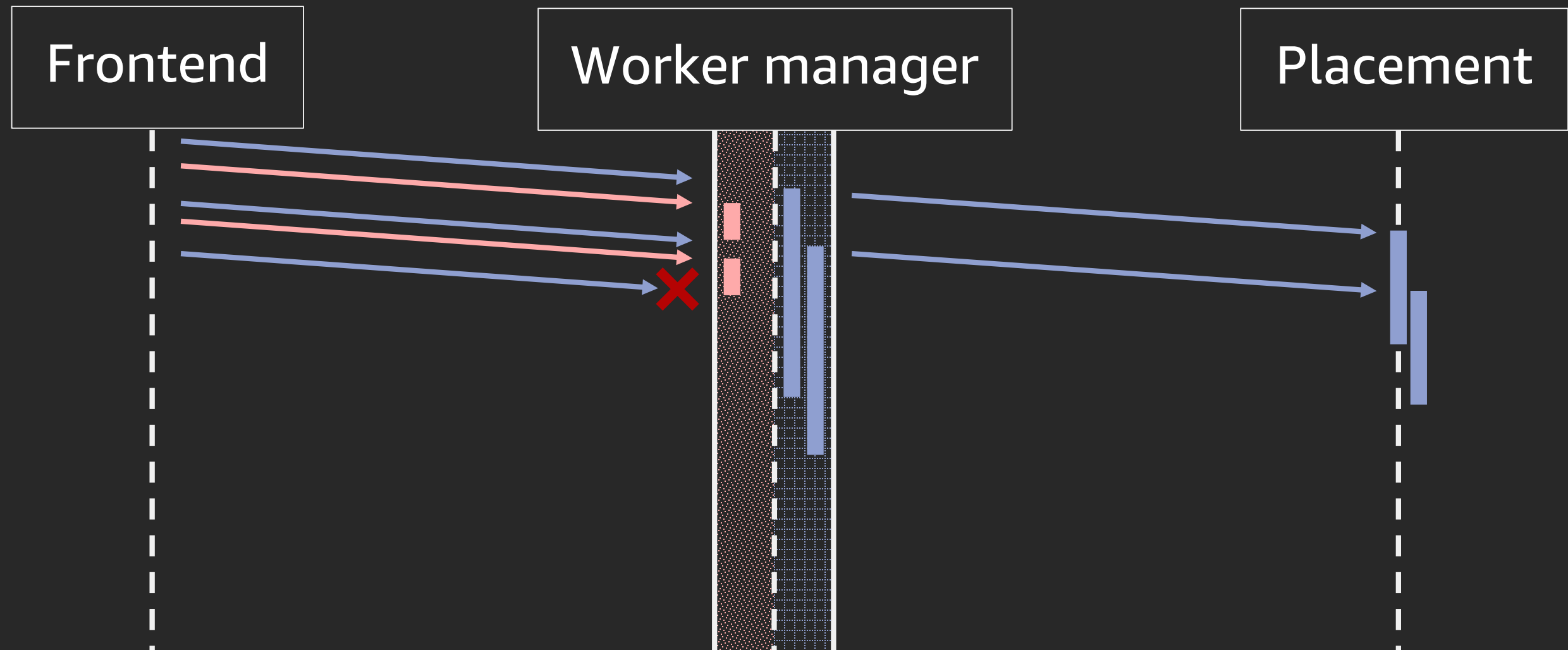**(only work on one thing at a time)**

# Bounded work

# Compartmentalize dependencies

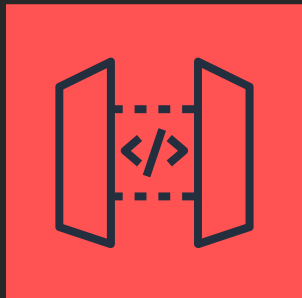# Compartmentalize dependencies

# Compartmentalize dependencies



AWS Lambda

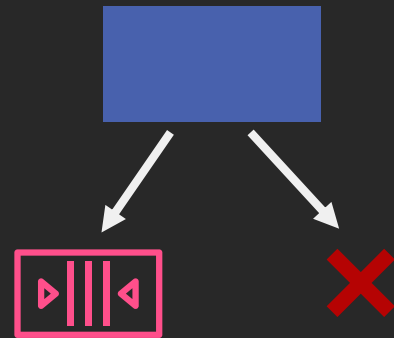(per-function concurrency control)



Amazon
API Gateway
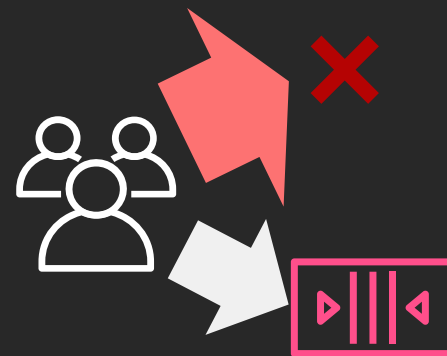
(per-API throttling)

# Watch out for queue backlogs

**Avoid pile-ups**

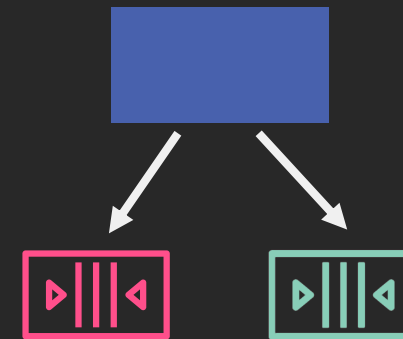Queues are quick to fill, slow to drain

**TTLs**

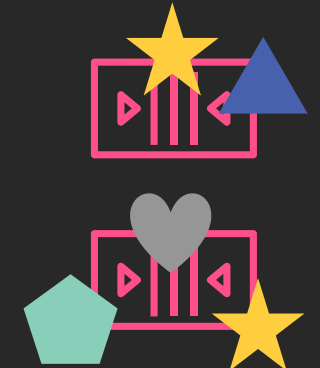Drop messages that are no longer relevant

**Backpressure**

Prevent a backlog from becoming unmanageable

**Priority queues**

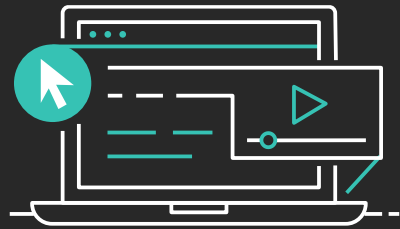Shift the backlog into a low-priority queue

**Shuffle-sharding**

Isolate backlogs in unrelated workloads
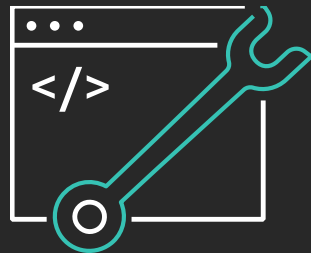
Serverless == resiliency

# Learn serverless with AWS Training and Certification

Resources created by the experts at AWS to help you learn modern application development

Free, on-demand courses on serverless, including

- Introduction to Serverless Development

- Getting into the Serverless Mindset

- AWS Lambda Foundations

- Amazon API Gateway for Serverless Applications

- Amazon DynamoDB for Serverless Architectures

Additional digital and classroom trainings cover modern application development and computing

Visit the Learning Library at https://aws.training

aws training and certification

# Thank you!

Please complete the session survey in the mobile app.