# Building confidence through chaos engineering on AWS

Naren Gakka (he/him)

Solutions Architect
AWS

# Agenda

(1) Chaos engineering (CE) introduction

(2) Continuous resilience (CR)

(3) Building a CR/CE program

(4) Chaos engineering resources

# Introduction to chaos engineering

# "Chaos engineering is about building a culture of resilience in the presence of unexpected system outcomes."
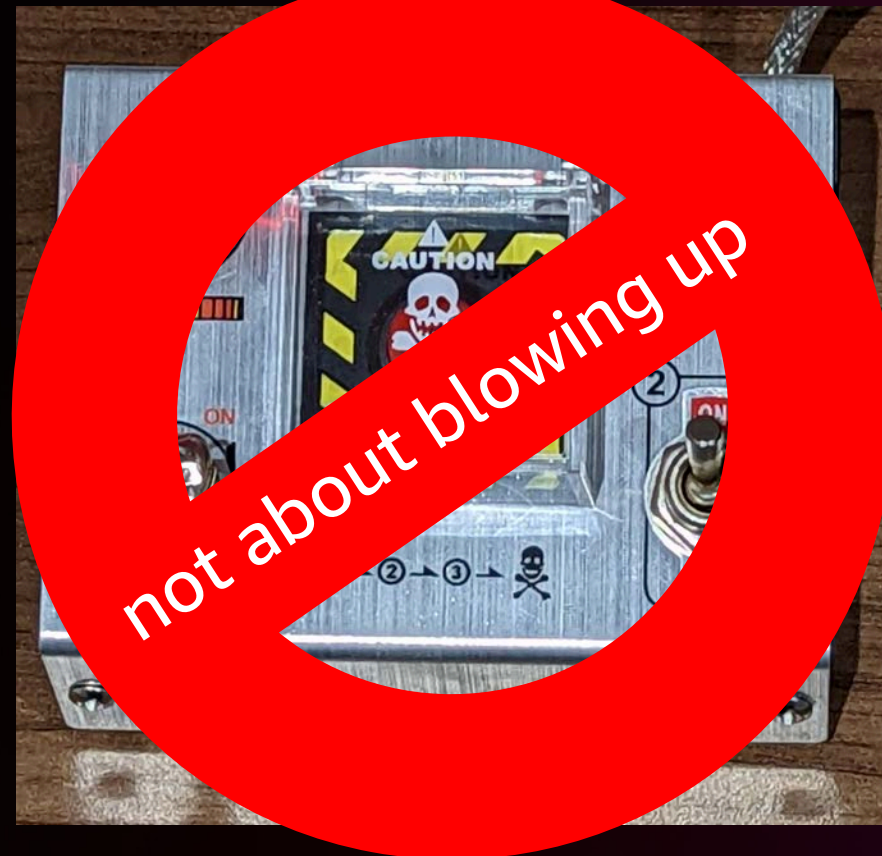
**Principles of chaos engineering**

principlesofchaos.org

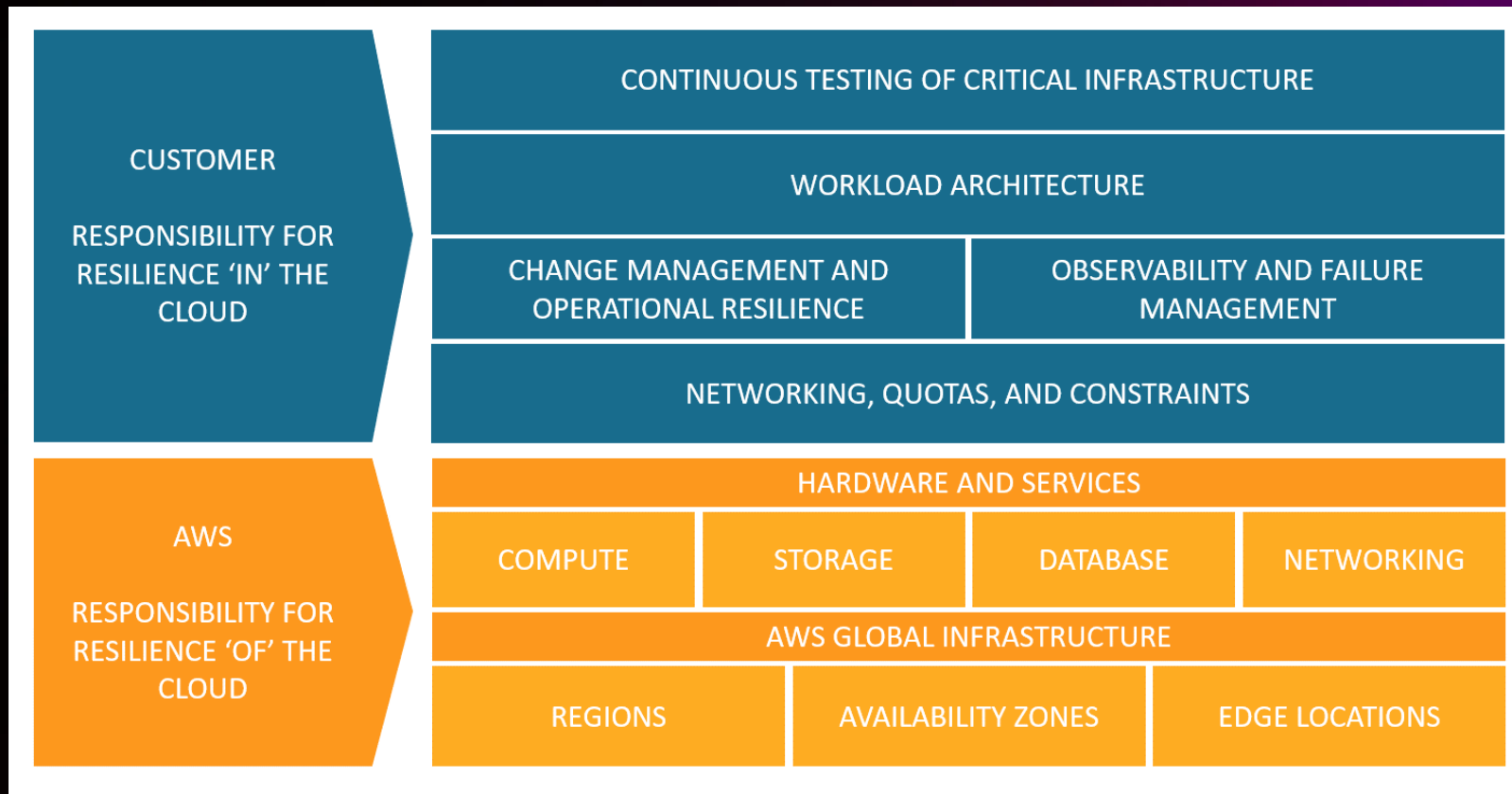# Chaos engineering a different perspective

Chaos engineering is  production

# Chaos engineering a different perspective

What you have control over

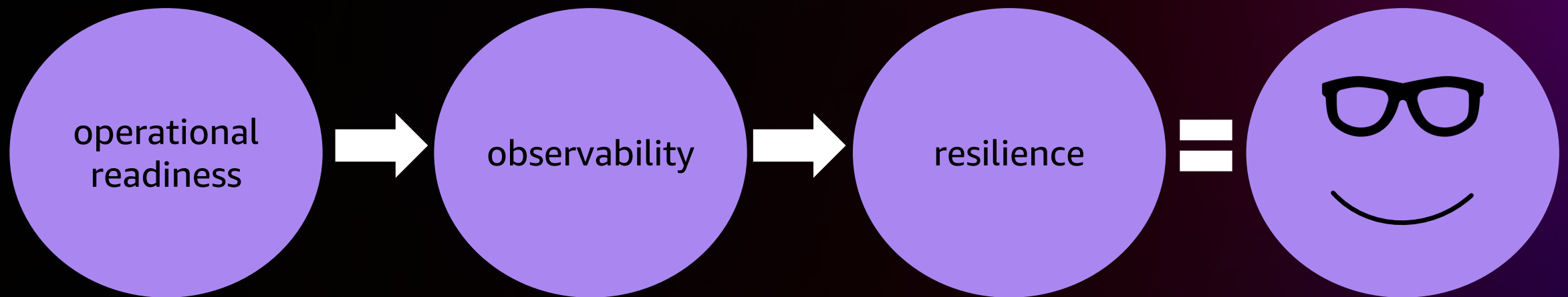What you have <span style="color:orange">no</span> control over

| CUSTOMER RESPONSIBILITY FOR RESILIENCE 'IN' THE CLOUD | CONTINUOUS TESTING OF CRITICAL INFRASTRUCTURE | |
|---|---|---|
| | WORKLOAD ARCHITECTURE | |
| | CHANGE MANAGEMENT AND OPERATIONAL RESILIENCE | OBSERVABILITY AND FAILURE MANAGEMENT |
| | NETWORKING, QUOTAS, AND CONSTRAINTS | |

| AWS RESPONSIBILITY FOR RESILIENCE 'OF' THE CLOUD | HARDWARE AND SERVICES | | | |
|---|---|---|---|---|
| | COMPUTE | STORAGE | DATABASE | NETWORKING |
| | AWS GLOBAL INFRASTRUCTURE | | | |
| | REGIONS | AVAILABILITY ZONES | EDGE LOCATIONS | |

**SHARED RESPONSIBILITY MODEL FOR RESILIENCE**
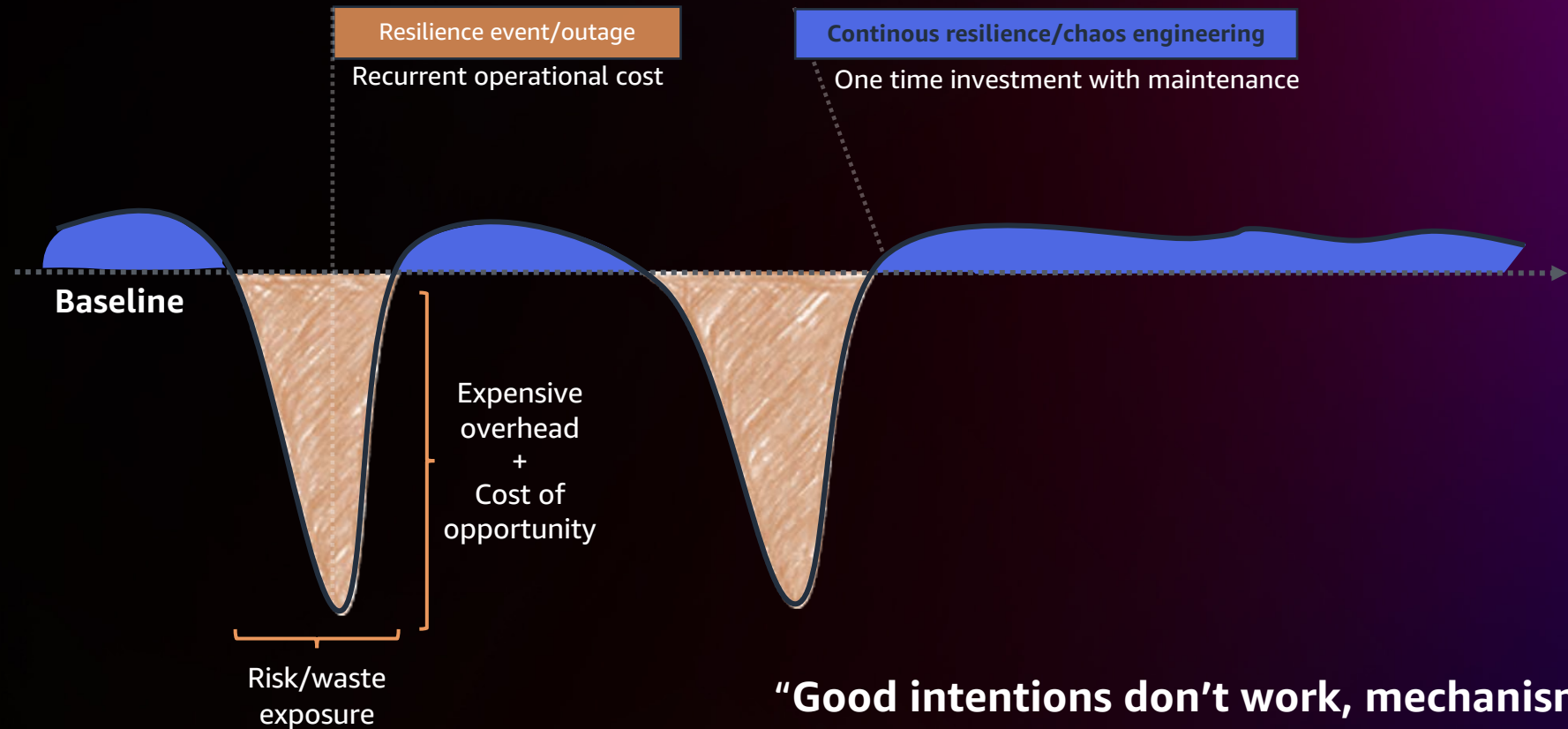
# Fail-safe chaos experiments

Perform controlled experiments in which the assumption is that the load or faults that you inject will be tolerated by the system and are fail-safe.

# How much confidence do you have in your system?

# Why chaos engineering?

Resilience

Technical debt

Security posture

Cost efficiency

Trade offs

Resilience event/outage

Recurrent operational cost

Continous resilience/chaos engineering

One time investment with maintenance

**Baseline**

Expensive
overhead
+
Cost of
opportunity

Risk/waste
exposure

**"Good intentions don't work, mechanisms do."**
– Jeff Bezos, founder of Amazon

# Chaos engineering stories

| Financial Services | Health Care | Insurance |
| --- | --- | --- |
| Media & Entertainment | Telco | Retail |
| Transport/Airlines | Hospitality | Food/Delivery |

https://github.com/ldomb/ChaosEngineeringPublicStories

# Chaos engineering adoption 2023

**40%** Will adopt chaos engineering as part  of their DevOps initiatives in 2023

Reducing unplanned downtime by **20%**

Source: The I&O Leader's Guide to Chaos Engineering by Gartner

# Chaos engineering prerequisites

**Prerequisite 1:**

Basic monitoring/observability

**Prerequisite 2:**

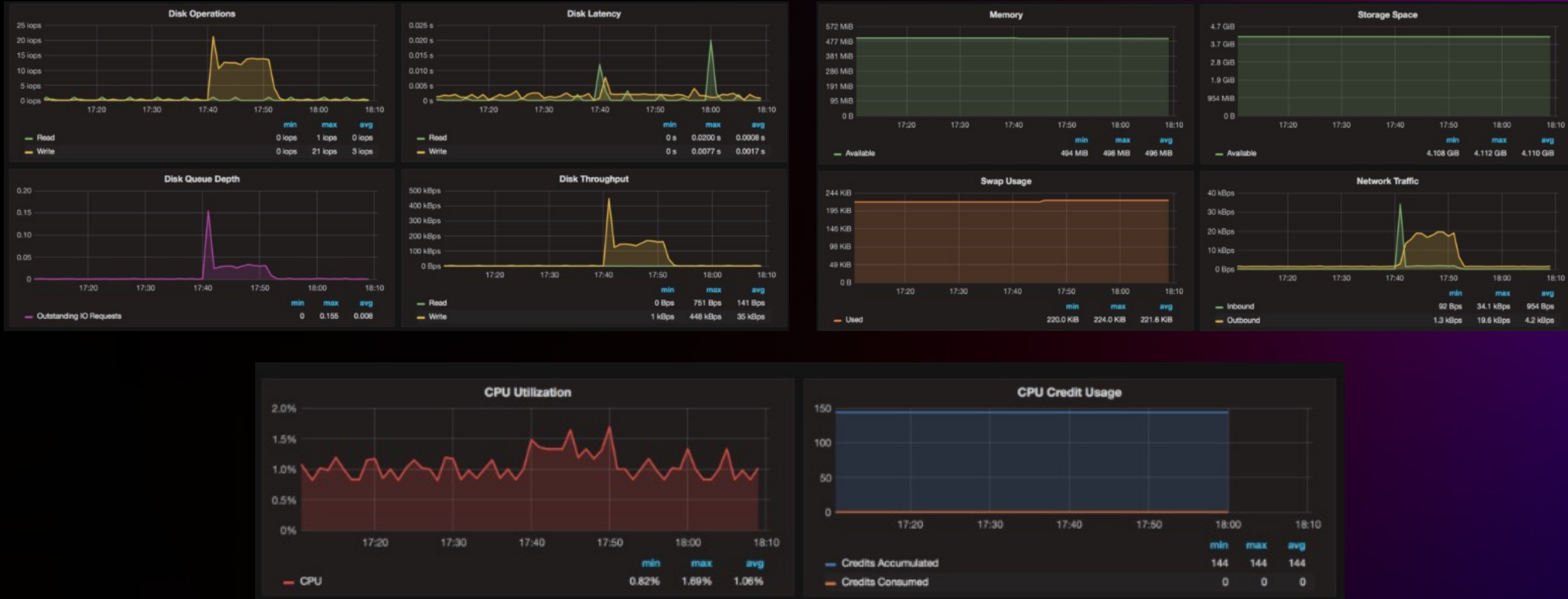Organizational awareness

**Prerequisite 3:**

Real world events/faults

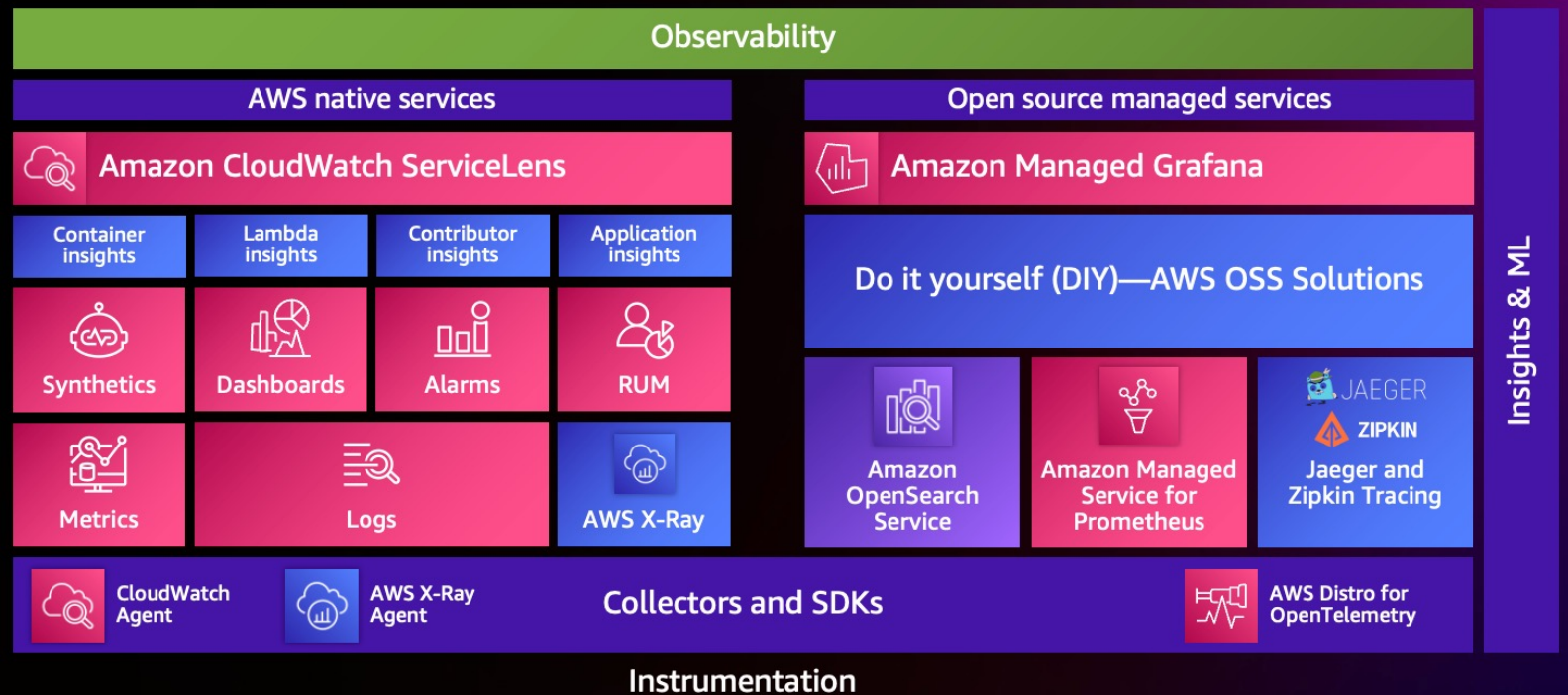**Prerequisite 4:**

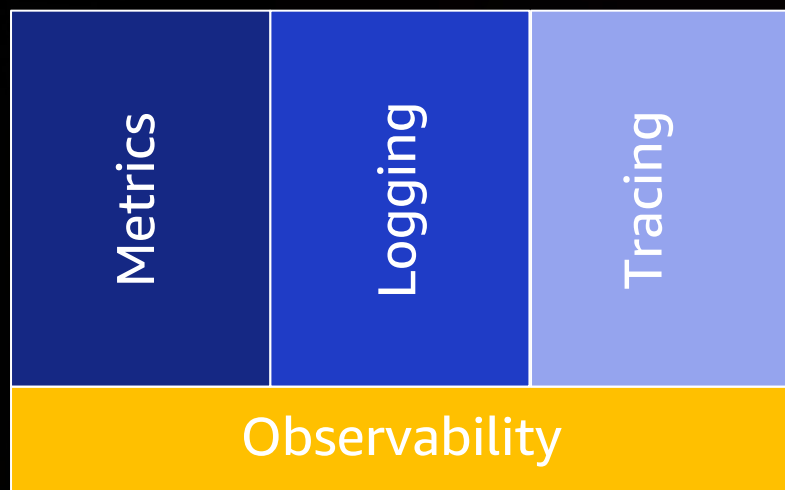Remediate findings

# Prerequisite 1: Basic monitoring/obsrevability

# Prerequisite 1: Basic monitoring/observability

## Observability

Find the needle in the haystack

# The 3 pillars of observability
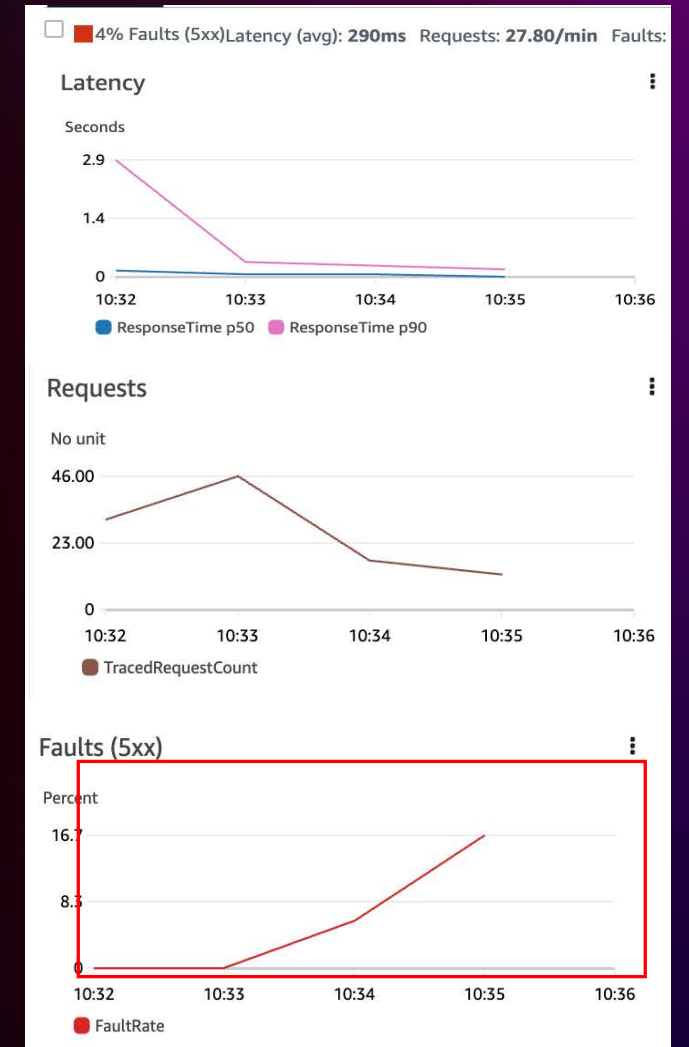
Metrics | Logging | Tracing

**Observability**
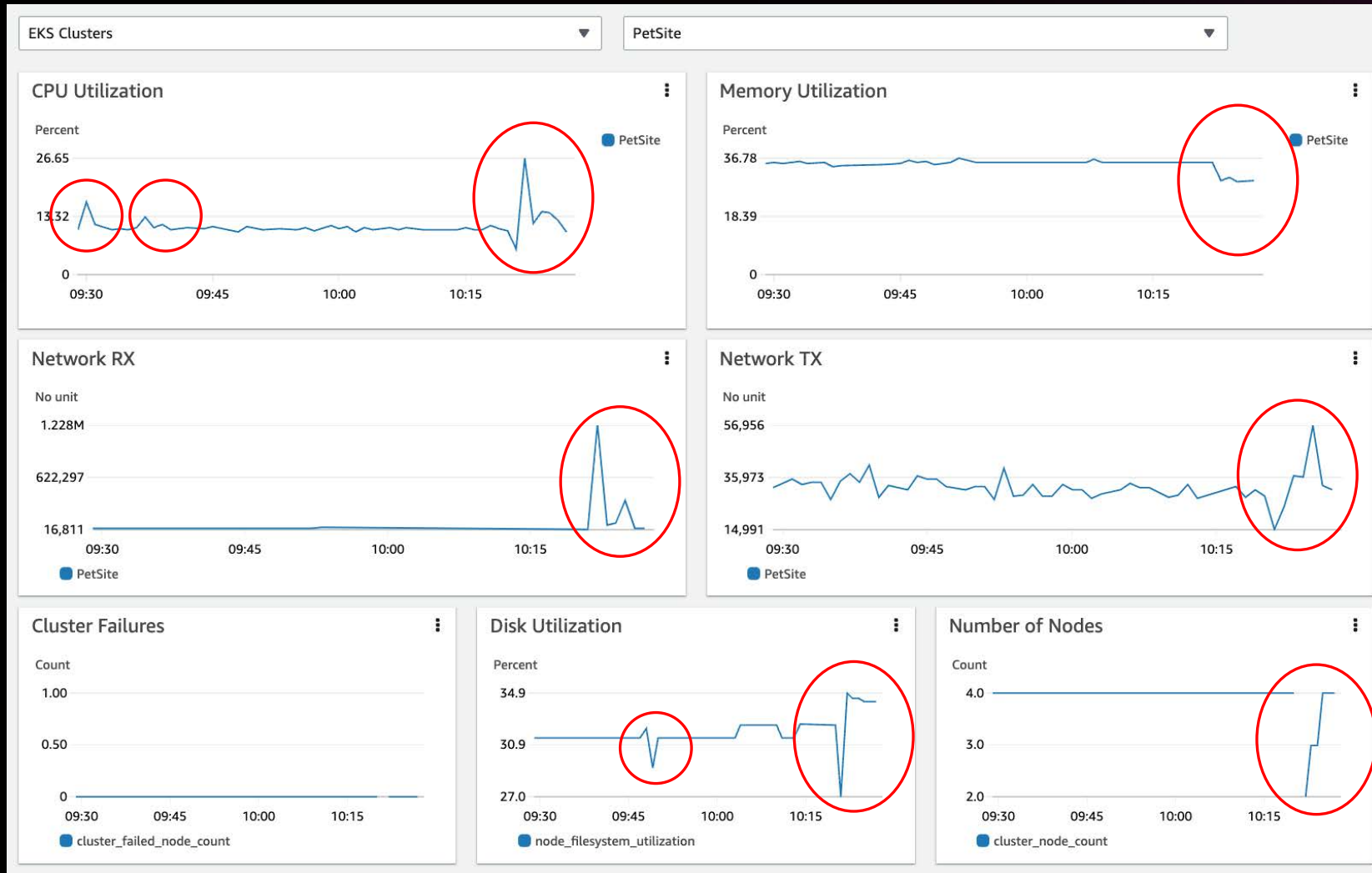
Do you understand the inner workings of your application?

Do you understand any system state your application may have fallen into?

Can you understand the above, just by observing your tools?

Can you understand the state, no matter how unusual the situation is?
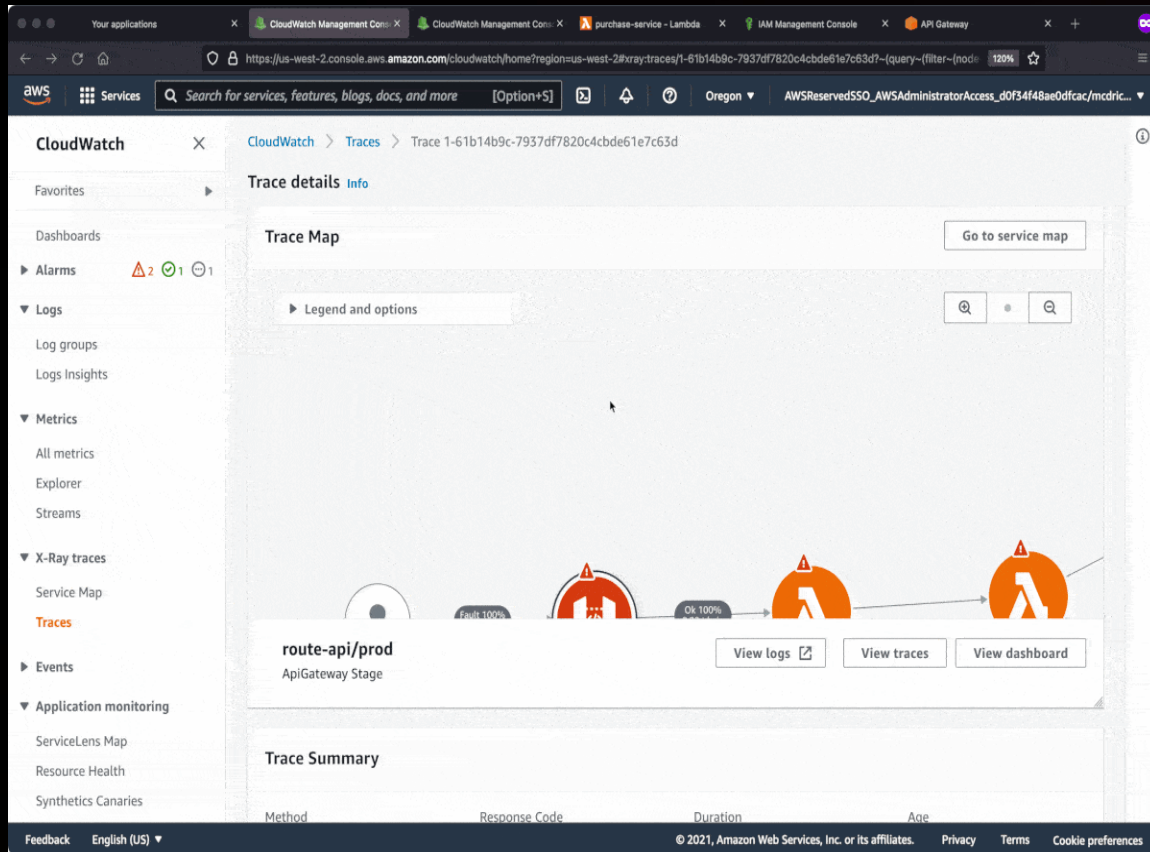
# Prerequisite 1: Basic monitoring/observability



Top level observability

# Prerequisite 1: Basic monitoring/observability



Application monitoring



User monitoring

https://aws.amazon.com/blogs/mt/an-observability-journey-with-amazon-cloudwatch-rum-evidently-and-servicelens/

# Prerequisite 2: Organizational awareness

**Central** team that **drives** chaos engineering

## Empowers

Education

Experiment Catalog

Learnings

🔴 Executive sponsor

Decentralized execution

Engineering team 1

- Experiment type
- Environment
- Self-service
- Guardrails
- Game Day
- Automated
- Resilience

Engineering team 2

- Experiment type
- Environment
- Self-service
- Guardrails
- Game Day
- Automated
- Resilience

# Prerequisite 3: Real-world experiments

**Code & configuration**
e.g., bad deployment, cred expiration, host shutdown

**Infrastructure**
e.g., datacenter failure, hardware failure

**Data and state**
e.g., data corruption, overload

**Dependencies**
e.g., third-party integrations, AWS services

**Highly unlikely, but technically feasible**
e.g., physical loss of an entire Region, the internet is down

# Prerequisite 4: Remediate the findings

- Findings through chaos engineering experiments should be prioritized based on the level of impact they may cause

- Findings that involve the resilience or security of your workload should have priority over new features, as if not addressed timely, can impact your customers

- Find an executive sponsor that can help you address the priority if needed

# Continuous resilience

# Continuous resilience

Code &
configuration
e.g. bad deployment,
cred expiration

Infrastructure
e.g. datacenter failure,
hardware failure

Data and state
e.g. data corruption

IMPROVE

STEADY
STATE

VERIFY

Gain
confidence in
your systems
capability to
withstand
component
failures

HYPOTHESIS

RUN
EXPERIMENT

Continuous resilience flywheel by Adrian Hornsby

Dependencies
e.g. third-party
integrations, AWS services

Highly unlikely, but
technically feasible
e.g. physical loss of an entire
Region, the internet is down

# Building a chaos engineering / continuous resilience program

# Setting expectations

1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

**Chaos experiment discovery**

Setting expectations → Select target workload → Validate workload readiness → Ensure observability/monitoring and playbooks → Address known-knowns → Define experiments; Scope Steady state; Hypothesis; Rollbacks → Experiment in lower level environment (canary) Verify → Run experiment in dev/staging → Establish guardrails in production → Experiment in production (canary) Verify → Run experiment in production → Continuous resilience (automation)

# Setting expectations

Project plan

Roles and responsibilities

Contribution

Outcome

# Selecting the target workload

PRIORITIZE FINDINGS, FIX THEM, RE-ITERATE. CONTINUOUSLY RUN EXPERIMENTS AS OFTEN AS WORKLOAD NEEDS; SCALE MECHANISM TO OTHER WORKLOADS.

1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

**Chaos experiment discovery**

Setting expectations →
Select target workload →
Validate workload readiness →
Ensure observability/ monitoring and playbooks →
Address known-knowns →
Define experiments; Scope Steady state; Hypothesis; Rollbacks →
Experiment in lower level environment (canary) Verify →
Run experiment in dev/staging →
Establish guardrails in production →
Experiment in production (canary) Verify →
Run experiment in production →
Continuous resilience (automation)

# Target workload selections

For your **first** experiment, choose a workload that if **degraded** has only **minimal** to **no** impact to your **internal** or **external** customers.

# Chaos experiment discovery

**PRIORITIZE FINDINGS, FIX THEM, RE-ITERATE. CONTINUOUSLY RUN EXPERIMENTS AS OFTEN AS WORKLOAD NEEDS; SCALE MECHANISM TO OTHER WORKLOADS.**

1 — What happened? (Timeline)
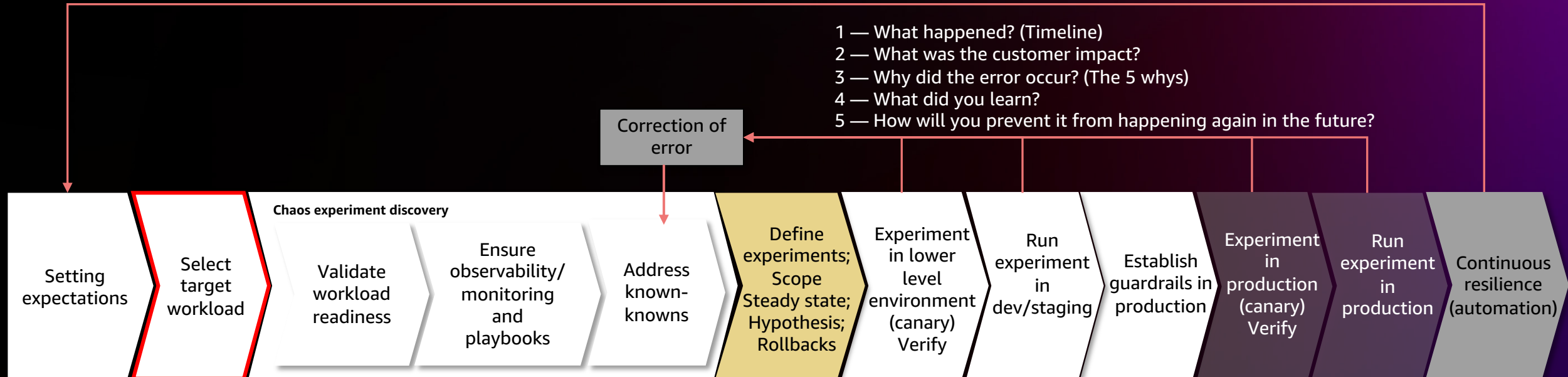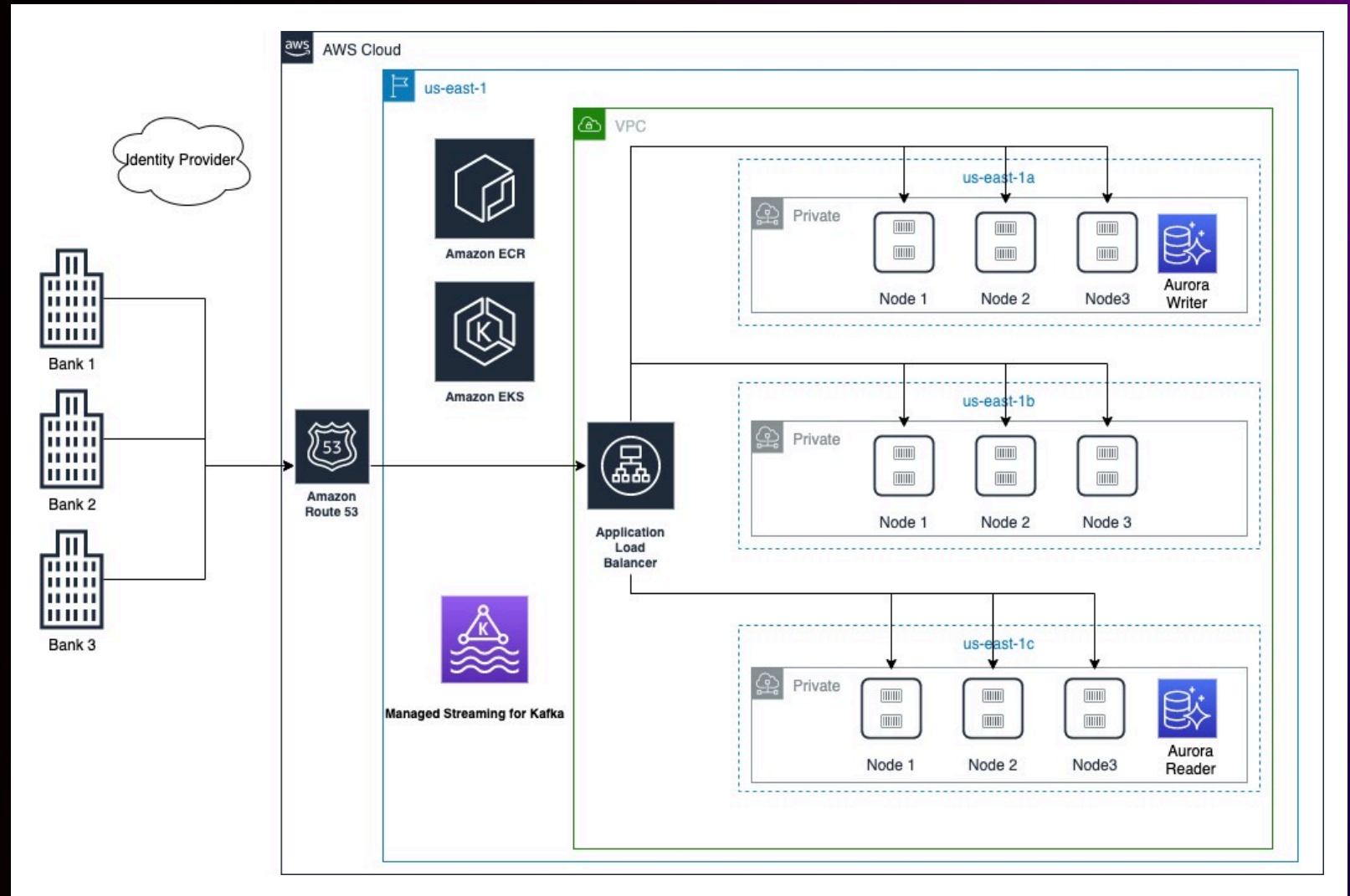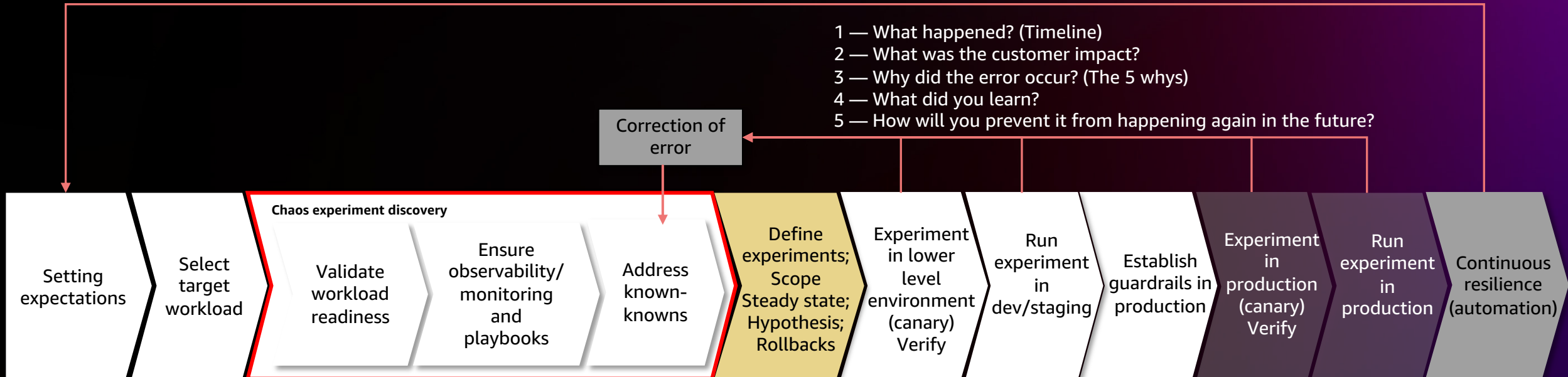2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

**Chaos experiment discovery**

- Setting expectations
- Select target workload
- Validate workload readiness
- Ensure observability/ monitoring and playbooks
- Address known-knowns
- Define experiments; Scope Steady state; Hypothesis; Rollbacks
- Experiment in lower level environment (canary) Verify
- Run experiment in dev/staging
- Establish guardrails in production
- Experiment in production (canary) Verify
- Run experiment in production
- Continuous resilience (automation)

# Chaos experiment discovery

Operational readiness review

Failover

Observability

Well-Architected Review

Runbooks

Retries

Incident response

Health checks

Circuit breakers

Fix known issues before moving forward with the experiment!

# Define experiment

1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

Setting expectations

Select target workload

**Chaos experiment discovery**

Validate workload readiness

Ensure observability/ monitoring and playbooks

Address known-knowns

Define experiments; Scope Steady state; Hypothesis; Rollbacks

Experiment in lower level environment (canary) Verify

Run experiment in dev/staging

Establish guardrails in production

Experiment in production (canary) Verify

Run experiment in production

Continuous resilience (automation)

# Chaos experiment definition – Failure modes

# Chaos experiment definition – Experiment #1

Brownout

Offered Throughput (TPS)

Response Time (ms)

200
180
160
140
120
100
80
60
40
20
0

# Chaos experiment definition – Steady state



Payments – Transactions per second

Retail – Order per second

Streaming – Stream starts per second

Media/Audio – Playback event started

# Chaos experiment definition - Hypothesis

At a rate of 300 TPS, if 40% of the nodes in the EKS node-group are terminated, the Transaction Create API continues to serve the 99th percentile of requests in under 100 ms (steady state)

The EKS nodes will recover within 5 minutes, and pods will get scheduled and process traffic within 8 minutes after the initiation of the experiment

Alerts will fire within 3 minutes

# Chaos experiment definition – Bring all together

## Chaos Experiment

**Realtime Payments**

Workload Name

**Resilience**

Contribution

**Brownout - Terminate 40% nodes**

Action

**Staging**

Environment

**30 minutes**

Duration

**300 TPS**

Load

**Amazon EKS Payments Cluster**

Targets

**Some clients calls will time out**

Fault isolation boundary

**CW Alarm when node count < 60%**

Stop Condition

**CFN template to built nodes**

Rollback

**Chaos-ready**

Resource Tag/ID/Filter

**OpenSearch/CWL/X-ray**
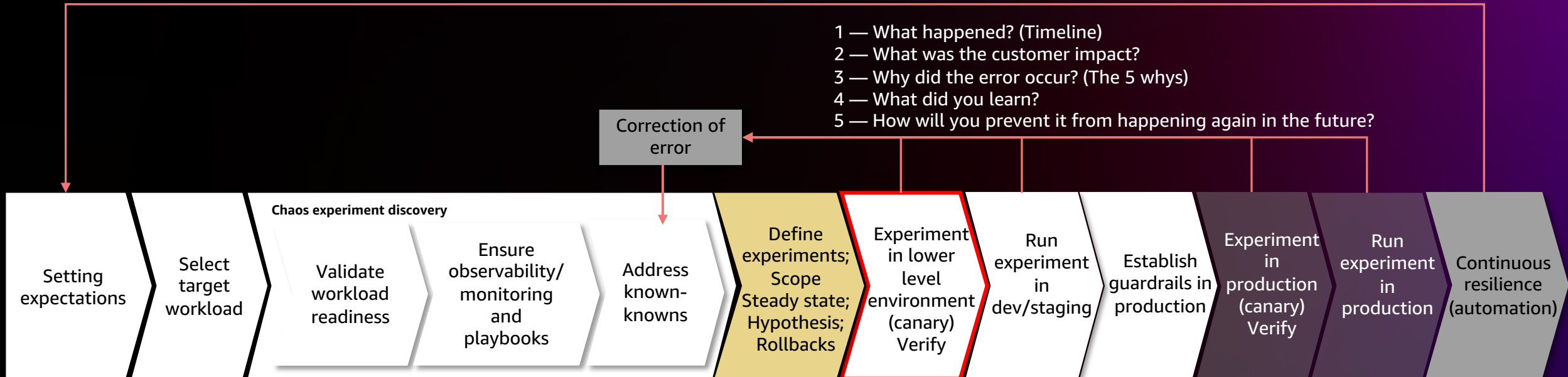
Observability / Logging

**At a rate of 300 TPS, If 40% …**

Hypothesis

Findings

COE – Actions to mitigate fault
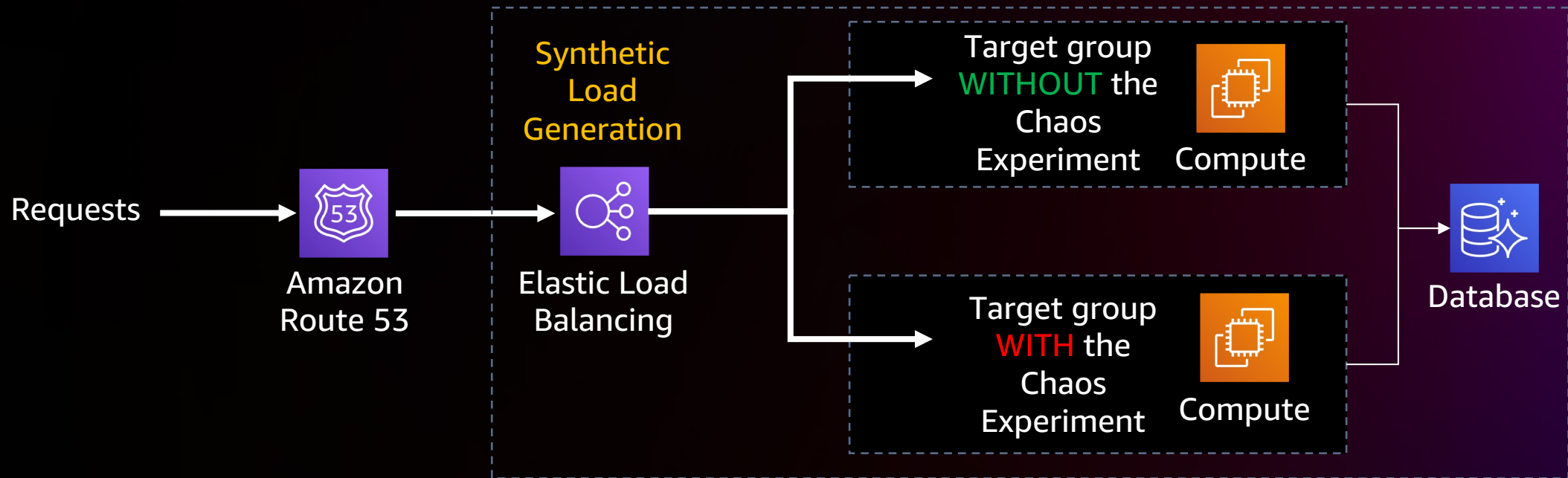
# Prime the environment for the experiment

1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

Chaos experiment discovery

Setting expectations

Select target workload

Validate workload readiness

Ensure observability/ monitoring and playbooks

Address known-knowns

Define experiments; Scope Steady state; Hypothesis; Rollbacks

Experiment in lower level environment (canary) Verify

Run experiment in dev/staging

Establish guardrails in production

Experiment in production (canary) Verify

Run experiment in production

Continuous resilience (automation)

# Chaos experiment execution flow



Is the system healthy? — No
Yes ↓

Is the experiment valid? — No
Yes ↓

Is control/experimental group defined? — No
Yes ↓

Generate load against target(s) — No
Yes ↓

Steady-state hypothesis within tolerance? — No
Yes ↓

Run action against target(s) — No
Yes ↓

Steady-state hypothesis within tolerance? — No → Stop condition/Roll back
Yes ↓

No deviations found

Deviations found

Experiment aborted

Fix and repeat

Repeat

# Controlled experiments through canary deployments in lower-level environment



Requests → Amazon Route 53 → **Synthetic Load Generation** → Elastic Load Balancing →

Target group **WITHOUT** the Chaos Experiment — Compute

Target group **WITH** the Chaos Experiment — Compute

→ Database

**Verify that both groups are healthy before moving forward**

# Run the experiment in development/staging

PRIORITIZE FINDINGS, FIX THEM, RE-ITERATE; CONTINUOUSLY RUN EXPERIMENTS AS OFTEN AS WORKLOAD NEEDS; SCALE MECHANISM TO OTHER WORKLOADS

1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

Chaos experiment discovery

Setting expectations

Select target workload

Validate workload readiness

Ensure observability/ monitoring and playbooks

Address known-knowns

Define experiments; Scope Steady state; Hypothesis; Rollbacks

Experiment in lower level environment (canary) Verify

Run experiment in dev/staging

Establish guardrails in production

Experiment in production (canary) Verify

Run experiment in production

Continuous resilience (automation)

# AWS Fault Injection
Actions

- Stop, reboot, and terminate instance(s) (Amazon EC2)

- API throttling/internal error/unavailable error (Amazon EC2)

- Increased memory or CPU load (Amazon EC2)

- Kill process (Amazon EC2)

- Latency injection (Amazon EC2)

- Drain container instances (Amazon ECS)

- Terminate task (Amazon ECS)

- Increase memory or CPU consumption per task (Amazon ECS)

- Terminate node group instances (Amazon EKS)

- Litmus Chaos/Chaos Mesh integration (Amazon EKS)

- Network Connectivity Disruption (Amazon EC2)

- Database stop, reboot, and failover (Amazon RDS)

# AWS Fault Injection
Actions

- Systems manager send command (SSM Agent)

- Increased memory or CPU load (Amazon EC2)

- Kill process (Amazon EC2)

- Latency injection (Amazon EC2)

- Increase memory or CPU consumption per task (Amazon EC2)

- Network port blackhole (Amazon EC2)

- Network latency

- Network latency at target source

- Network packet loss

- Network packet loss at target source

# Execute the controlled experiment

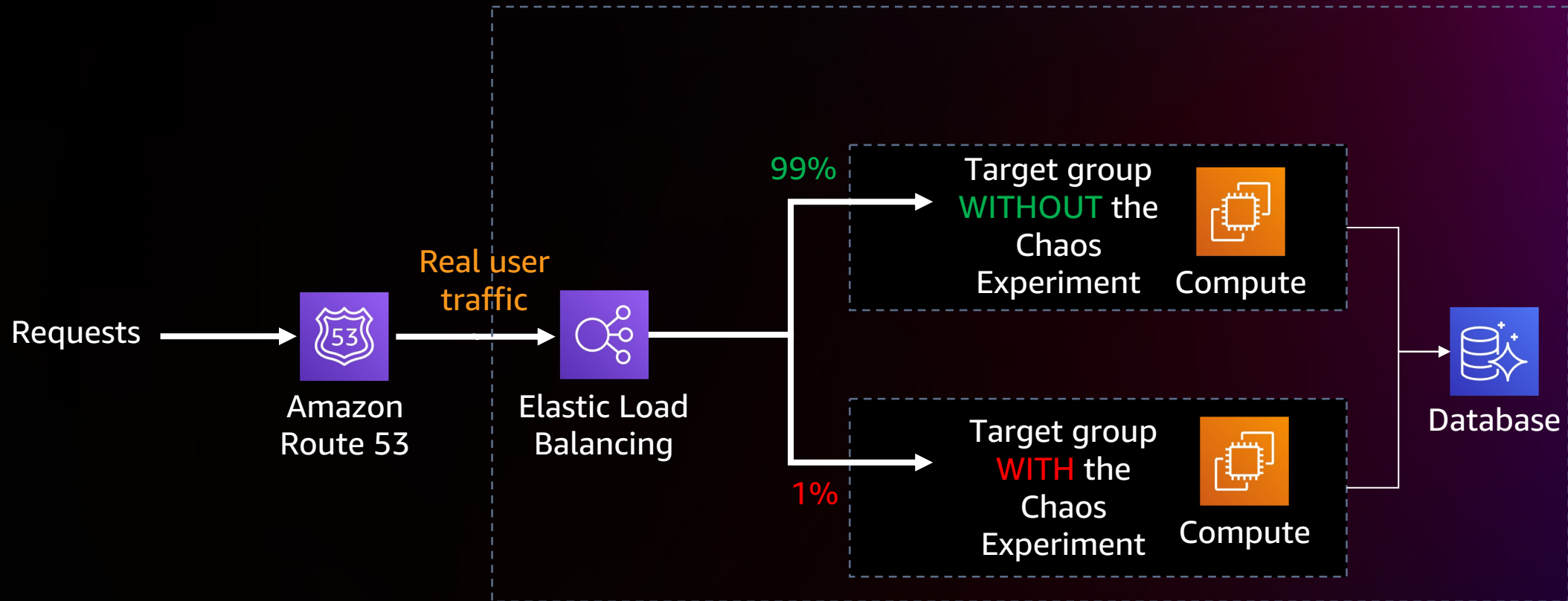## AWS Fault Injection Simulator

# Establish guardrails in production

PRIORITIZE FINDINGS, FIX THEM, RE-ITERATE; CONTINUOUSLY RUN EXPERIMENTS AS OFTEN AS WORKLOAD NEEDS; SCALE MECHANISM TO OTHER WORKLOADS
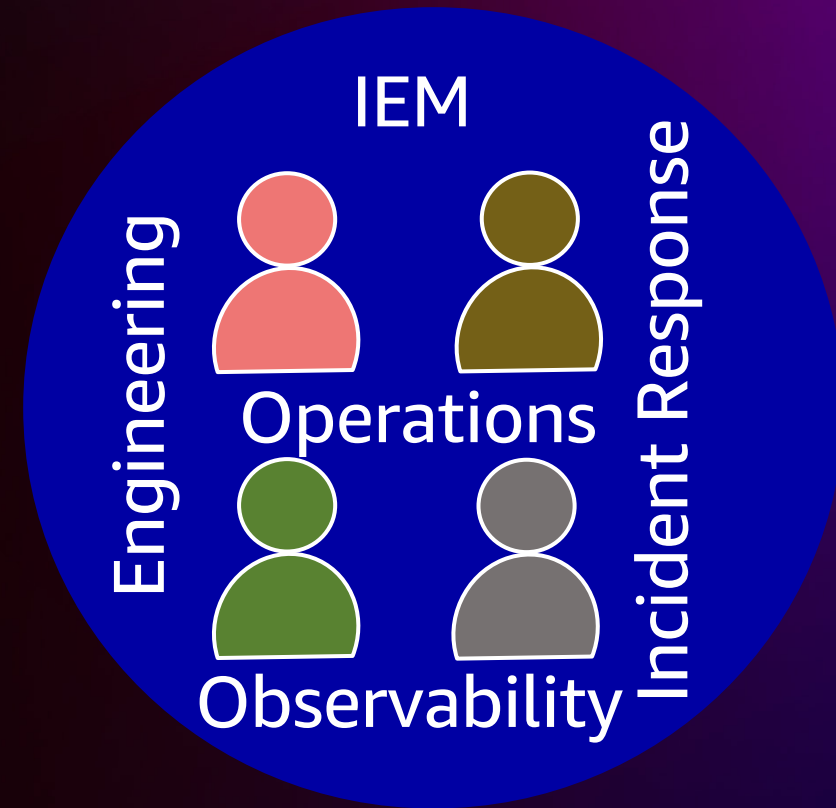
1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

**Chaos experiment discovery**

Setting expectations → Select target workload → Validate workload readiness → Ensure observability/ monitoring and playbooks → Address known-knowns → Define experiments; Scope Steady state; Hypothesis; Rollbacks → Experiment in lower level environment (canary) Verify → Run experiment in dev/staging → Establish guardrails in production → Experiment in production (canary) Verify → Run experiment in production → Continuous resilience (automation)

# Establish guardrails in production

- Run experiments off peak hours

- Validate that your IAM Permissions are sufficient for your experiment in production

- Validate if your fault isolation boundary is the same as in the lower-level environment or changes in the production environment

- Decide if synthetic traffic will be generated, or if you are planning to run the experiment against a subset of your customers

- Validate that observability is in place

- Validate that your runbooks/playbook are up to date in production

# Prime the prod environment for the experiment
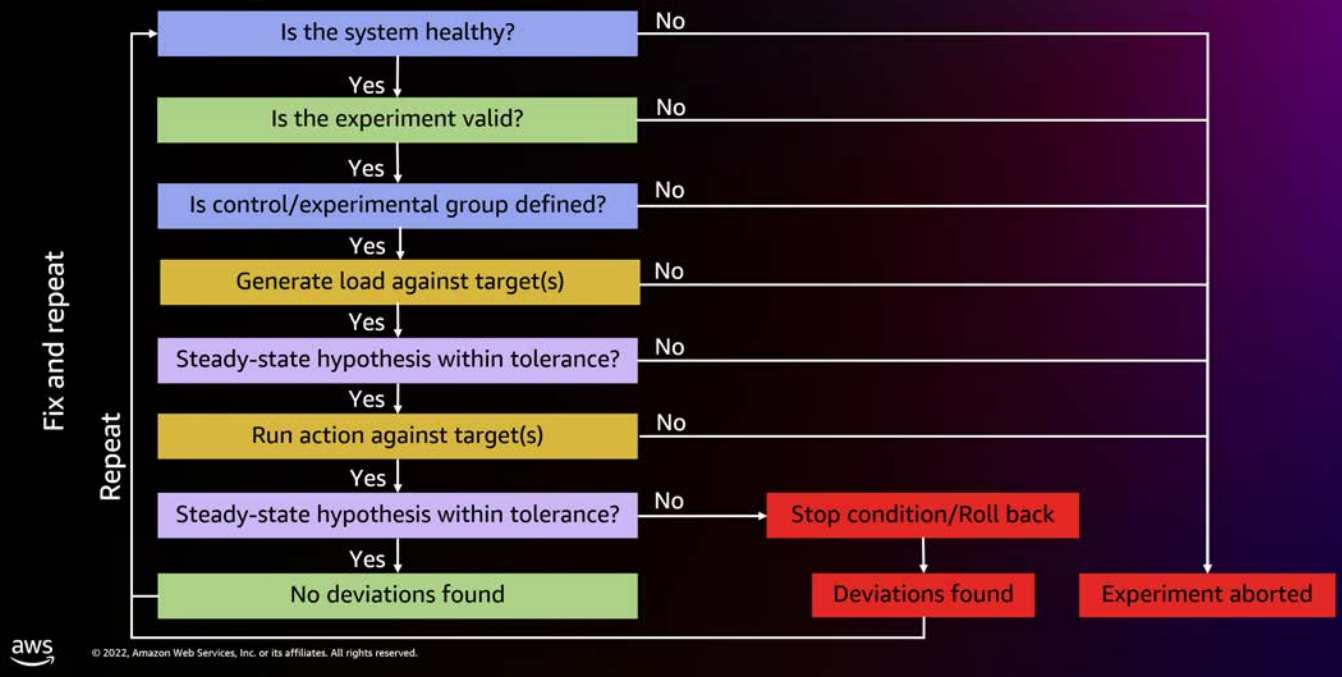
PRIORITIZE FINDINGS, FIX THEM, RE-ITERATE; CONTINUOUSLY RUN EXPERIMENTS AS OFTEN AS WORKLOAD NEEDS; SCALE MECHANISM TO OTHER WORKLOADS
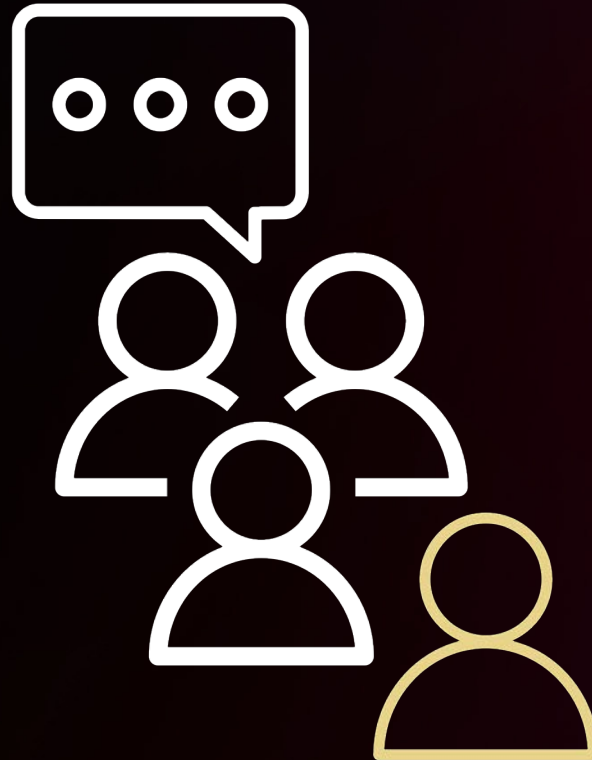
1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

**Chaos experiment discovery**

| Setting expectations | Select target workload | Validate workload readiness | Ensure observability/ monitoring and playbooks | Address known-knowns | Define experiments; Scope Steady state; Hypothesis; Rollbacks | Experiment in lower level environment (canary) Verify | Run experiment in dev/staging | Establish guardrails in production | Experiment in production (canary) Verify | Run experiment in production | Continuous resilience (automation) |

# Controlled experiments through canary deployments

# Controlled experiments through canary deployments

Requests → Amazon Route 53 → **Real user traffic** → Elastic Load Balancing

**99%** → Target group WITHOUT the Chaos Experiment — Compute

**1%** → Target group WITH the Chaos Experiment — Compute

→ Database

# Run the experiment in production

PRIORITIZE FINDINGS, FIX THEM, RE-ITERATE; CONTINUOUSLY RUN EXPERIMENTS AS OFTEN AS WORKLOAD NEEDS; SCALE MECHANISM TO OTHER WORKLOADS

1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

**Chaos experiment discovery**

Setting expectations

Select target workload

Validate workload readiness

Ensure observability/ monitoring and playbooks

Address known-knowns

Define experiments; Scope Steady state; Hypothesis; Rollbacks

Experiment in lower level environment (canary) Verify

Run experiment in dev/staging

Establish guardrails in production

Experiment in production (canary) Verify

Run experiment in production

Continuous resilience (automation)

# Execute the controlled experiment in production

# Run the experiment in production

PRIORITIZE FINDINGS, FIX THEM, RE-ITERATE; CONTINUOUSLY RUN EXPERIMENTS AS OFTEN AS WORKLOAD NEEDS; SCALE MECHANISM TO OTHER WORKLOADS

1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

Setting expectations

Select target workload

**Chaos experiment discovery**

Validate workload readiness

Ensure observability/ monitoring and playbooks

Address known-knowns

Define experiments; Scope Steady state; Hypothesis; Rollbacks

Experiment in lower level environment (canary) Verify

Run experiment in dev/staging

Establish guardrails in production

Experiment in production (canary) Verify

Run experiment in production

Continuous resilience (automation)

# Post mortem – Correction of error

How did we communicate

What did we learn?

Was everyone needed present?

How do we share what we've learned?

# Continuous resilience

PRIORITIZE FINDINGS, FIX THEM, REITERATE; CONTINUOUSLY RUN EXPERIMENTS AS OFTEN AS WORKLOAD NEEDS; SCALE MECHANISM TO OTHER WORKLOADS

1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

**Chaos experiment discovery**

| Setting expectations | Select target workload | Validate workload readiness | Ensure observability/ monitoring and playbooks | Address known-knowns | Define experiments; Scope Steady state; Hypothesis; Rollbacks | Experiment in lower level environment (canary) Verify | Run experiment in dev/staging | Establish guardrails in production | Experiment in production (canary) Verify | Run experiment in production | Continuous resilience (automation) |

# Automate the experiments

AUTOMATE

Individual experiments

AWS GameDays

Scheduled experiments

# Supporting documents

1 — What happened? (Timeline)
2 — What was the customer impact?
3 — Why did the error occur? (The 5 whys)
4 — What did you learn?
5 — How will you prevent it from happening again in the future?

Correction of error

**Chaos experiment design phase**

| Setting expectations | Select target workload | Validate workload readiness | Ensure observability/ monitoring and playbooks | Address known-knowns | Define experiments; Scope Steady state; Hypothesis; Rollbacks | Experiment in lower level environment (canary) Verify | Run experiment in dev/staging | Establish guardrails in production | Experiment in production (canary) Verify | Run experiment in production | Continuous resilience (automation) |

| Introduce the program | Process maturity discovery | Chaos engineering handbook (CEH) | Chaos experiment template | Guardrails review assessment | Infrastructure event management (TAM) | CR Implementation Guidance |
| Introduction to chaos engineering | Workload maturity discovery | Well-architected review | Interim findings report | | Final report: Next steps and findings | |
| | | | Chaos experiment review | | | |

Process journey map defining steps and stakeholders (RACI)

# Chaos engineering
# AWS resources

# Chaos engineering/observability workshops

# Chaos engineering/observability workshops


Chaos engineering on AWS


Validating security guardrails with chaos engineering


Resilience engineering


Serverless chaos workshop


Observability workshop


AWS fault isolation boundaries


Multi-AZ resilience patterns


Chaos engineering stories

# Thank You!

Naren Gakka
twitter: @narengka