AWS
re:Invent

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

COP343

# Observability best practices at Amazon

David Yanacek (he/him)

Sr. Principal Engineer
AWS Monitoring & Observability

Ian McGarry (he/him)

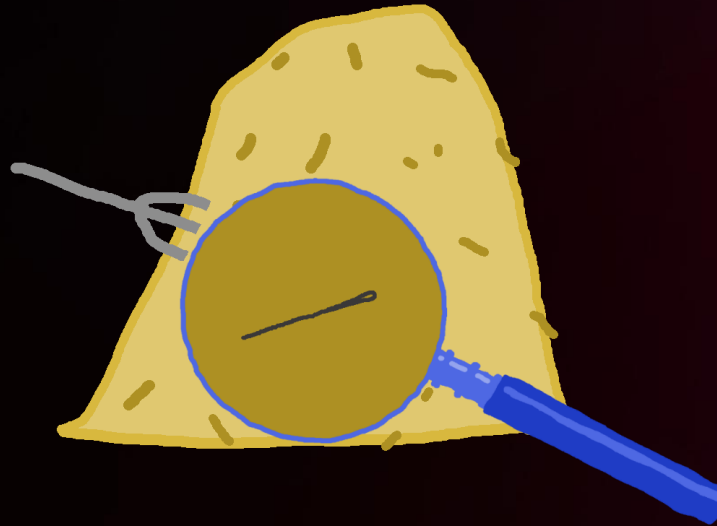Director of Engineering
AWS Monitoring & Observability

# Table of contents



The DevOps
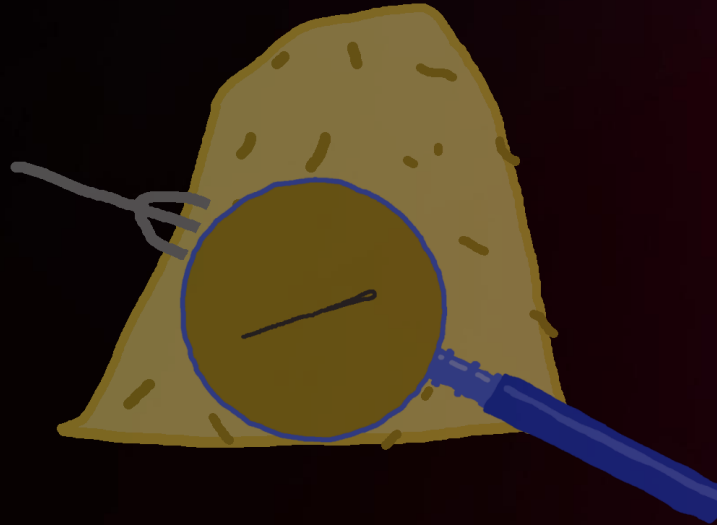flywheel at Amazon



Find the
root cause



Measure from
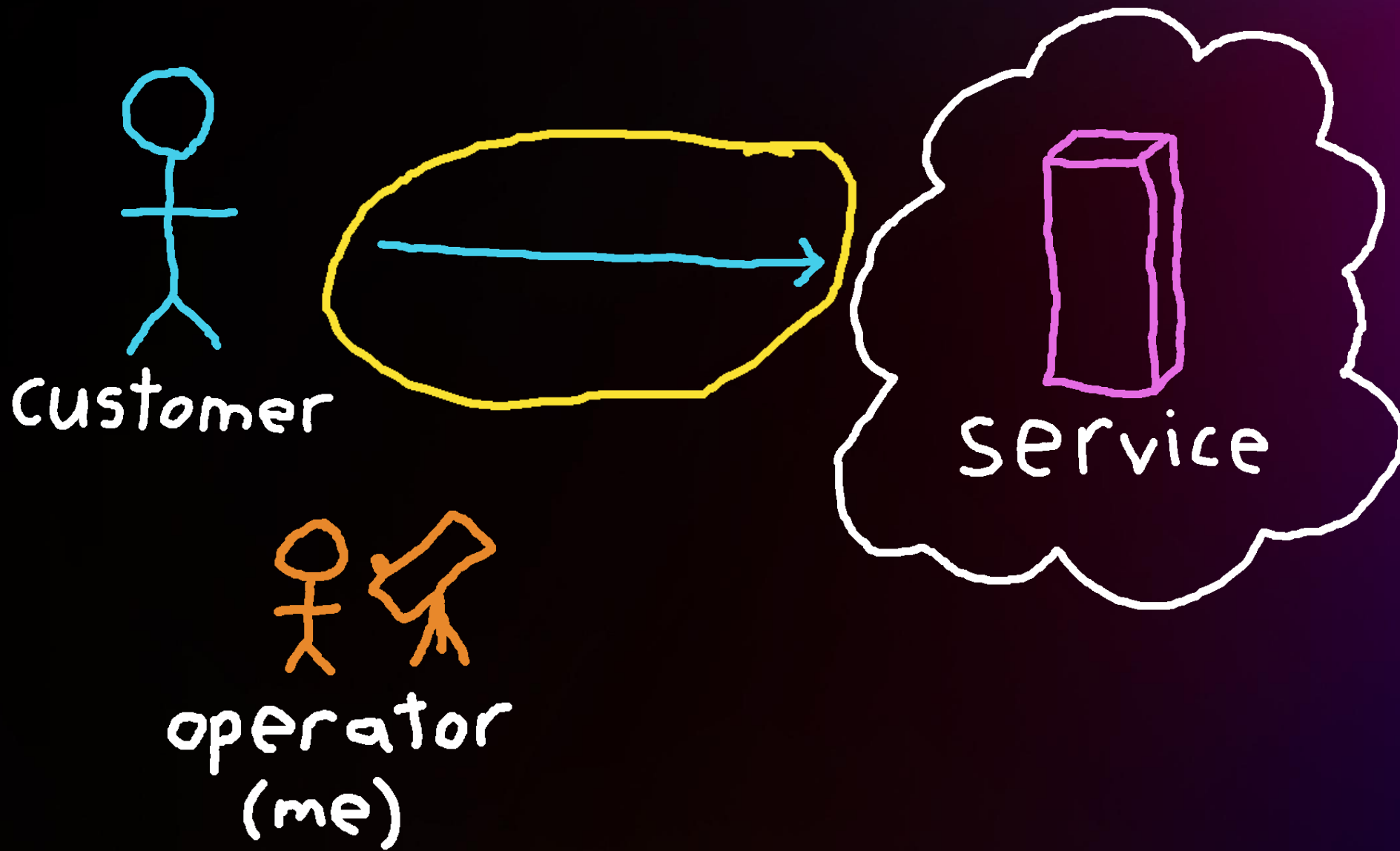everywhere

# Chapter one

The DevOps
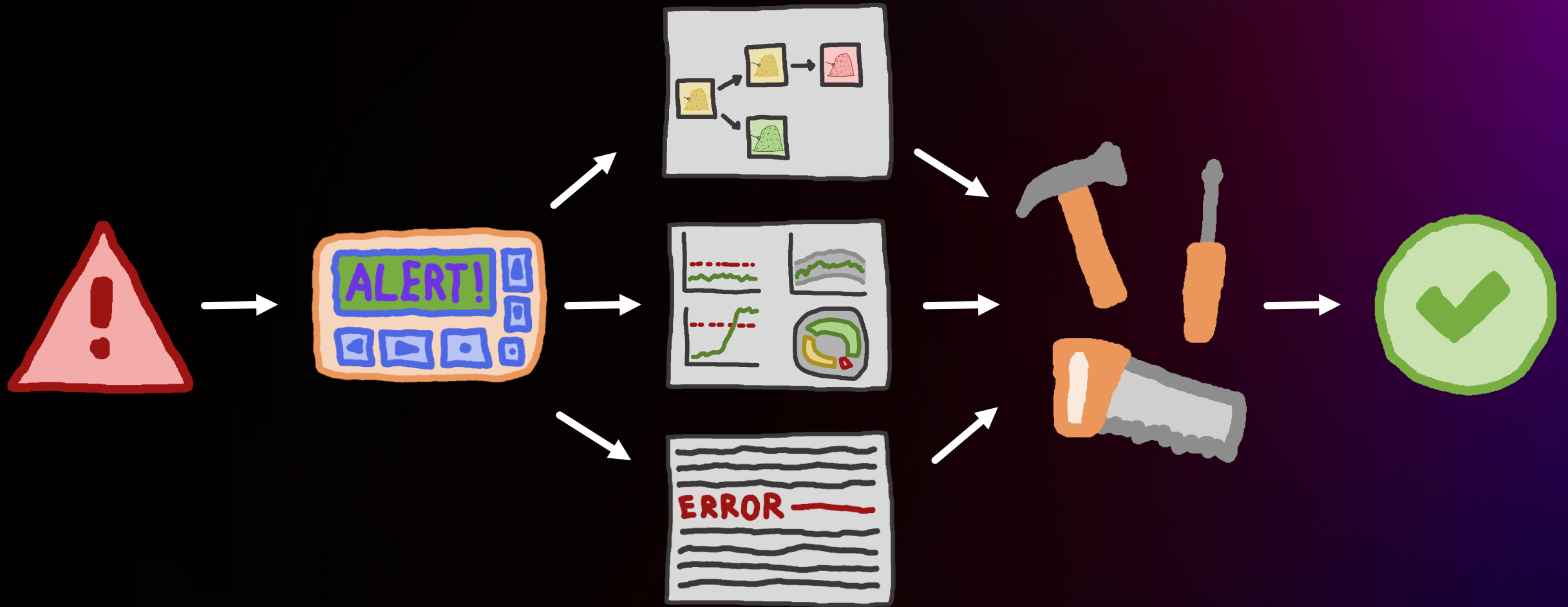flywheel at Amazon

Find the
root cause

Measure from
everywhere

# Observing customer experience

# Responding to incidents

# Post-incident retrospectives

# Weekly operations meeting(s)

# Observability at scale

*"Amazon CloudWatch is used to monitor more than 9 quadrillion metric observations, ingests more than 5 exabytes of logs per month."* *As of November 2022*

# The cycle of observability

# Chapter two

The DevOps
flywheel at Amazon

**Find the
root cause**

Measure from
everywhere

Finding the needle in the haystack

Alarm: Awareness of a needle

Which haystack?

Dashboard: Find the right haystack

Traces: Following the clues

Traces: Following the clues

Metrics: Find the part of the haystack

# Logs: Find the needle

# Finding the root cause

- Propagate trace information
- Build lots of dashboards
- Get multi-dimensional
- Use high cardinality metrics
- Dig even deeper with log analysis
- Look at the raw logs
- Analyze using profiler

# Finding the root cause

- **Propagate trace information**
- Build lots of dashboards
- Get multi-dimensional
- Use high cardinality metrics
- Dig even deeper with log analysis
- Look at the raw logs
- Analyze using profiler

Finding the right haystack with traces

# Reasoning about distributed systems

# Distributed tracing

Trace-Id: 1-5759e988

Users → Application Load Balancer → Frontend → Async queue ← Poller

# Distributed system-wide visibility



Users → Application Load Balancer → Frontend → Async queue ← Poller

# Finding the root cause

- Propagate trace information
- **Build lots of dashboards**
- Get multi-dimensional
- Use high cardinality metrics
- Dig even deeper with log analysis
- Look at the raw logs
- Analyze using profiler

# Hierarchy of dashboards

# Customer experience dashboards

# System-level dashboards

# Microservice dashboards

# Capacity dashboards

# Hierarchy of dashboards



Customer experience dashboard

Clients

System-level dashboard

aws

API Microservice

Instances

Microservice dashboard

Backend Microservice

Containers

Database

Backend Microservice

Lambdas

Capacity dashboard

# Dashboard design

# The dashboard is never done

# Finding the root cause

- Propagate trace information
- Build lots of dashboards
- **Get multi-dimensional**
- Use high cardinality metrics
- Dig even deeper with log analysis
- Look at the raw logs
- Analyze using profiler

# Use metric dimensions to drill down

# Log everything: Embedded Metrics Format

```json
{
  "_aws": {
    "Timestamp": 1574109732004,
    "CloudWatchMetrics": [
      {
        "Namespace": "ProductInfoService",
        "Dimensions": [],
        "Metrics": [
          { "Name": "Time", "Unit": "Milliseconds" }
        ]
      }
    ],
  },
  "CustomerId": "user@example.com",
  "InstanceId": "us-west-2c",
  "Time": 100,
  "TraceId": "Root=1-5759e988-bd862e3fe1be46a994272793;Sampled=1"
}
```

Instrumentation

Configuration

# Example: Simple web service code

```
public ProductInfo getProductInfo(String customerId, String productId) {

    ProductInfo info = cache.get(customerId, productId);

    if (info == null) {
        info = db.query(customerId, productId);
    }

    return info;
}
```

What product or customer was this for?

Was this a cache hit or miss?

If it failed, was it a timeout or did it return an error?

How long did the db query take? Did it return anything?

# Example: Instrumenting code

```
public ProductInfo getProductInfo(String customerId, String productId,
    Metrics metrics) {

    metrics.addProperty("CustomerId", customerId);
    metrics.addProperty("ProductId", productId);


    ProductInfo info = cache.get(customerId, productId, metrics);


    if (info == null) {
        info = db.query(customerId, productId, metrics);
    }


    return info;
}
```

Record what we're working on

Propagate per-request
context object

# Instrumentation in the log

```json
{
...
  "CustomerId": "user@example.com",
  "ProductId": "BrightOrangeSneakers",
  "ClientIp": "192.168.131.39",
  "InstanceId": "i-001234a4bf70dec41EXAMPLE",
  "AvailabiltyZone": "us-west-2a",
  "Cache.Node": "10.0.2.182",
  "Cache.Hit": 0,
  "Cache.Time": 1,
  "DB.Time": 10,
  "DB.Results": 1,
  "Time": 100,
  "TraceId": "Root=1-5759e988-bd862e3fe1be46a994272793;Sampled=1"
}
```

Caller and request info

Infrastructure info

Timing and info around dependencies

# Log everything about a unit of work

```json
{
  ...
  "CustomerId": "user@example.com",
  "ProductId": "BrightOrangeSneakers",
  "ClientIp": "192.168.131.39",
  "InstanceId": "i-001234a4bf70dec41EXAMPLE",
  "AvailabiltyZone": "us-west-2a",
  "Cache.Node": "10.0.2.182",
  "Cache.Hit": 0,
  "Cache.Time": 1,
  "DB.Time": 10,
  "DB.Results": 1,
  "Time": 100,
  "TraceId": "Root=1-5759e988-bd862e3fe1be46a994272793;Sampled=1"
}
```

Attributes

Measurements

# Log everything: Embedded Metrics Format

```json
{
  "_aws": {
    "Timestamp": 1574109732004,
    "CloudWatchMetrics": [
      {
        "Namespace": "ProductInfoService",
        "Dimensions": [],
        "Metrics": [
          { "Name": "Time", "Unit": "Milliseconds" }
    ]}]},
  "CustomerId": "user@example.com",
  "AvailabilityZone": "us-west-2c",
  "Time": 100,
  "TraceId": "Root=1-5759e988-bd862e3fe1be46a994272793;Sampled=1"
}
```

Metric definitions

Defines which data to
turn into a metric

ProductInfoService Latency

1h  3h  12h  1d  3d  1w  Custom  Line  Actions  1m

Milliseconds

Time

140.00
120.00
100.00
80.00
60.00
40.00
20.00
0

20:15   20:20   20:25   20:30   21:15

```
{
    "Namespace": "ProductInfoService",
    "Dimensions": [],
    "Metrics": [
        { "Name": "Time", "Unit": "Milliseconds" }
    ]
}

    "Time": 103,
```

Browse   Query   Graphed metrics (1)   Options                Add query

**Metrics** (1) Info                                          Graph search

Oregon ▼   All > ProductInfoService > Metrics with no dimensions   Search for any metric, dimension or resource id   < 1 >

ProductInfoService ✎ ✕

☑ Metric name 1/1                                             ▲

☑ Time

# ProductInfoService Latency

| 1h | 3h | 12h | 1d | 3d | 1w | Custom ▦ | Line ▾ | Actions ▾ | ⟳ | 1m ▾ |

Milliseconds

■ Time



| 140.00 |
| 120.00 |
| 100.00 |
| 80.00 |
| 60.00 |
| 40.00 |
| 20.00 |
| 0 |

20:15  20:20  20:25  20:30  20:35  20:40  20:45  20:50  20:55  21:00  21:05  21:10  21:15

=

| Browse | Query | Graphed metrics (1) | Options | Source |

Add math ▾    Add query ▾

## Metrics (1) Info

Graph with SQL    Graph search

| Oregon ▾ | All > ProductInfoService > Metrics with no dimensions | 🔍 Search for any metric, dimension or resource id | ‹ 1 › |

ProductInfoService ☑ ✕                                                          ⚙

| ☑ | **Metric name 1/1** | ▲ |
| ☑ | Time |  |

# Product Info Service architecture



Users → Application Load Balancer → Product Info Service → Search index, Update queue, Product DB

# Product Info Service architecture

GetProduct()

Users → Application Load Balancer → Product Info Service → Search index / Update queue / Product DB

# Product Info Service architecture

UpdateProduct()



Users

Application
Load Balancer

Product
Info Service

Search index

Update queue

Product DB

# Product Info Service architecture

SearchProducts()



Users → Application Load Balancer → Product Info Service → Search index

Product Info Service → Update queue

Product Info Service → Product DB

# Finding the root cause

- Propagate trace information
- Build lots of dashboards
- Get multi-dimensional
- **Use high cardinality metrics**
- Dig even deeper with log analysis
- Look at the raw logs
- Analyze using profiler

PER-PRODUCT

# Dimensionality

"Show me **[metric]** per **[dimension]**"

| Latency |
| Errors |
| **Measurements** |

| **Availability Zone** |
| API |
| EC2 Instance |
| Product Id |
| **Customer Id** |
| **Attributes** |

(Few of these)

Cardinality

(Many of these)

# Analyzing high cardinality metrics

"Are all of my Amazon EC2 instances returning errors?"

"Are different customer experience different impact?"

"Is just one Product Id having problems?"


LATENCY PER-INSTANCE

# CloudWatch Metric Insights

```
SELECT SUM(Errors)
FROM ProductInfoService
GROUP BY Operation
ORDER BY SUM() DESC
LIMIT 5
```



1 - GetProductAttributes
2 - QueryProduct
3 - UpdateProductPrice

CloudWatch  >  Metrics

Errors by API ✎

Errors

3,011

1,506

0

17:00          18:00          19:00

1 - GetProductReviews
2 - UpdateProductAttribute
3 - ListProductsByRating
4 - ListProductsBySubcategory
5 - GetProductImages

...

(many more)

# CloudWatch Contributor Insights

# Combining dimensionality and cardinality

# Client vs. server fault

## Client fault (4XX)

## Server fault (5XX)

```
+------+--------------------------+
| 500  | Internal Server Error    |
| 501  | Not Implemented          |
| 502  | Bad Gateway              |
| 503  | Service Unavailable      |
| 504  | Gateway Timeout          |
| 505  | HTTP Version Not Supported |
+------+--------------------------+
```

```
+------+------------------------------+
| 400  | Bad Request                  |
| 401  | Unauthorized                 |
| 402  | Payment Required             |
| 403  | Forbidden                    |
| 404  | Not Found                    |
| 405  | Method Not Allowed           |
| 406  | Not Acceptable               |
| 407  | Proxy Authentication Required |
| 408  | Request Timeout              |
| 409  | Conflict                     |
| 410  | Gone                         |
| 411  | Length Required              |
| 412  | Precondition Failed          |
| 413  | Payload Too Large            |
| 414  | URI Too Long                 |
| 415  | Unsupported Media Type       |
| 416  | Range Not Satisfiable        |
| 417  | Expectation Failed           |
| 426  | Upgrade Required             |
+------+------------------------------+
```

# Input validation bugs

(Bug: Making input validation **more** strict)

```
    public void validate() {
-       if(input.length() > 140) {
+       if(input.length() > 100) {
            throw new ClientError();
        }
    }
```

# Graphing top contributors

# Graphing top contributors

# Metrics from another dimension

# Dimensionality and cardinality together

$$\frac{\text{\# of customers with errors}}{\text{\# of customers}}$$

$$\frac{INSIGHT\_RULE\_METRIC(ErrorsPerAccount, UniqueContributors)}{INSIGHT\_RULE\_METRIC(RequestsPerAccount, UniqueContributors)}$$
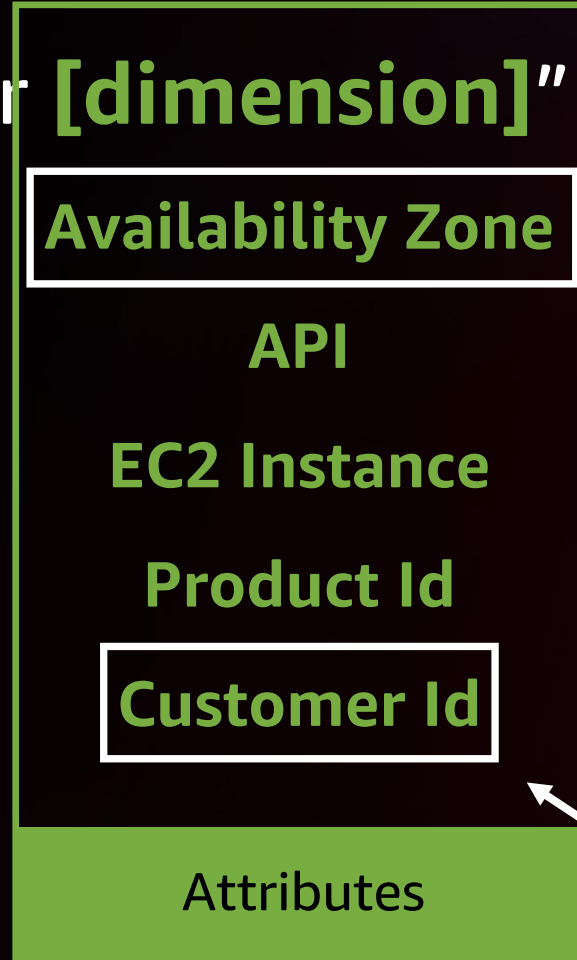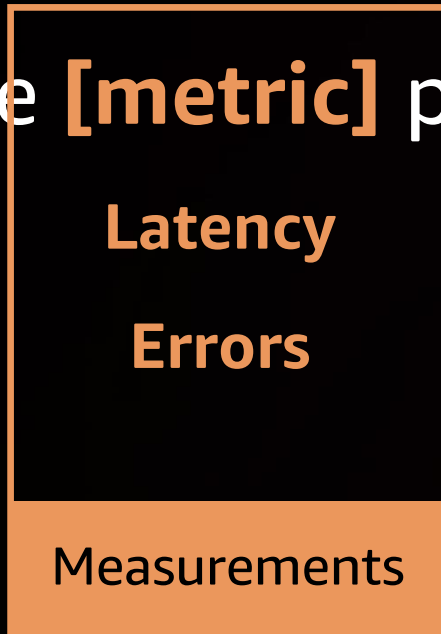
# Metrics from another dimension

# Finding the root cause

- Propagate trace information
- Build lots of dashboards
- Get multi-dimensional
- Use high cardinality metrics
- **Dig even deeper with log analysis**
- Look at the raw logs
- Analyze using profiler

SELECT
FROM
WHERE

# IoT monitoring application



AWS IoT

Analytics app

# IoT monitoring application

# When request rates are a lie

# IoT monitoring application

# IoT monitoring application

# IoT monitoring application



AWS IoT

14:30:01

Analytics app

# Slice and dice log analysis



```
1  fields @timestamp, @message
2    | sort @timestamp
3    | stats count(@message) by bin(1s)
4
```

# Aggregation defined up front

```json
{
  "Operation": "GetProduct",
  "CustomerId": "user@example.com",
  "ProductId": "BrightOrangeSneakers",
  "ClientIp": "192.168.131.39",
  "InstanceId": "i-001234a4bf70dec41EXAMPLE",
  "AvailabiltyZone": "us-west-2a",
  "Cache.Node": "10.0.2.182",
  "Cache.Hit": 0,
  "Cache.Time": 1,
  "DB.Time": 10,
  "Error": 0,
  "Time": 100,
  "TraceId": "Root=1-5759e988-bd862e3fe1be46a994272793;Sampled=1"
}
```

# Aggregation defined up front

```
{
  "Operation": "GetProduct",
  "CustomerId": "user@example.com",
  "ProductId": "BrightOrangeSneakers",
  "ClientIp": "192.168.131.39",
  "InstanceId": "i-001234a4bf70dec
  "AvailabiltyZone": "us-west-2a",
  "Cache.Node": "10.0.2.182",
  "Cache.Hit": 0,
  "Cache.Time": 1,
  "DB.Time": 10,
  "Error": 0,
  "Time": 100,
  "TraceId": "Root=1-5759e988-bd862e3fe1be46a994272793;Sampled=1"
}
```

"Show me latency per-operation"

```
"Dimensions": [["Operation"]],
"Metrics": [
  { "Name": "Time",
    "Unit": "Milliseconds" }]
```

# One-off aggregation

```json
{
    "Operation": "GetProduct",
    "CustomerId": "user@example.com",
    "ProductId": "BrightOrangeSneakers",
    "ClientIp": "192.168.131.39",
    "InstanceId": "i-001234a4bf70ded
    "AvailabiltyZone": "us-west-2a",
    "Cache.Node": "10.0.2.182",
    "Cache.Hit": 0,
    "Cache.Time": 1,
    "DB.Time": 10,
    "Error": 0,
    "Time": 100,
    "TraceId": "Root=1-5759e988-bd862e3fe1be46a994272793;Sampled=1"
}
```

"Show me the top 20 products that are returning errors for requests from IP address 192.168.131.39"

```
1  filter ClientIp='192.168.131.39' and Error > 0
2  | stats count() as errorCount by ProductId
3  | sort errorCount desc
4  | limit 20
```

# One-off aggregation

```json
{
    "Operation": "GetProduct",
    "CustomerId": "user@example.com",
    "ProductId": "BrightOrangeSneakers",
    "ClientIp": "192.168.131.39",
    "InstanceId": "i-001234a4bf70ded",
    "AvailabiltyZone": "us-west-2a",
    "Cache.Node": "10.0.2.182",
    "Cache.Hit": 0,
    "Cache.Time": 1,
    "DB.Time": 10,
    "Error": 0,
    "Time": 100,
    "TraceId": "Root=1-5759e988-bd862e3fe1be46a994272793;Sampled=1"
}
```

"Show me the products with the lowest cache hit rate, and which cache node they landed on

```
1  stats (sum(Cache.Hit) / count()) as cacheHitRate
2    by ProductId, Cache.Node
3  | sort cacheHitRate
4  | limit 20
```
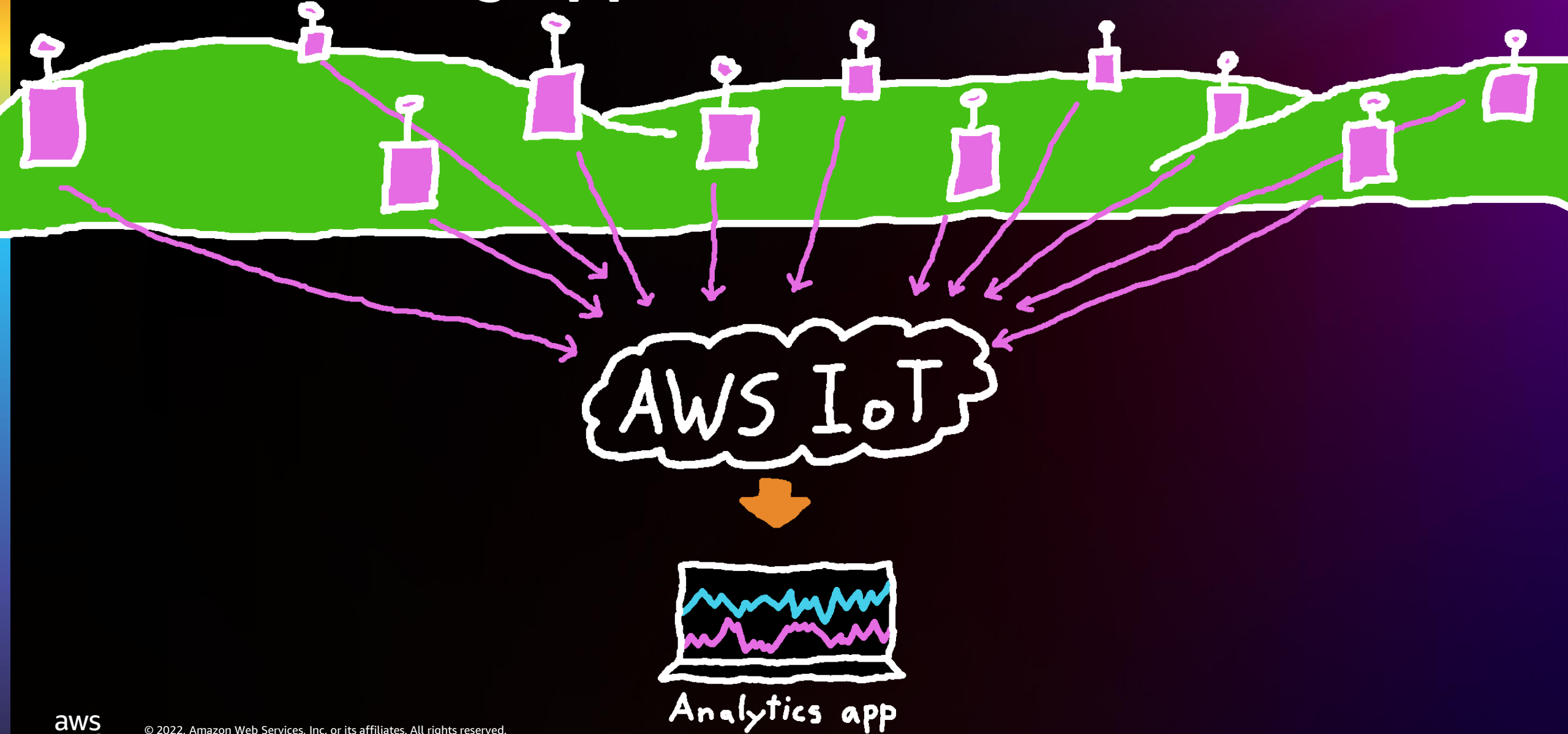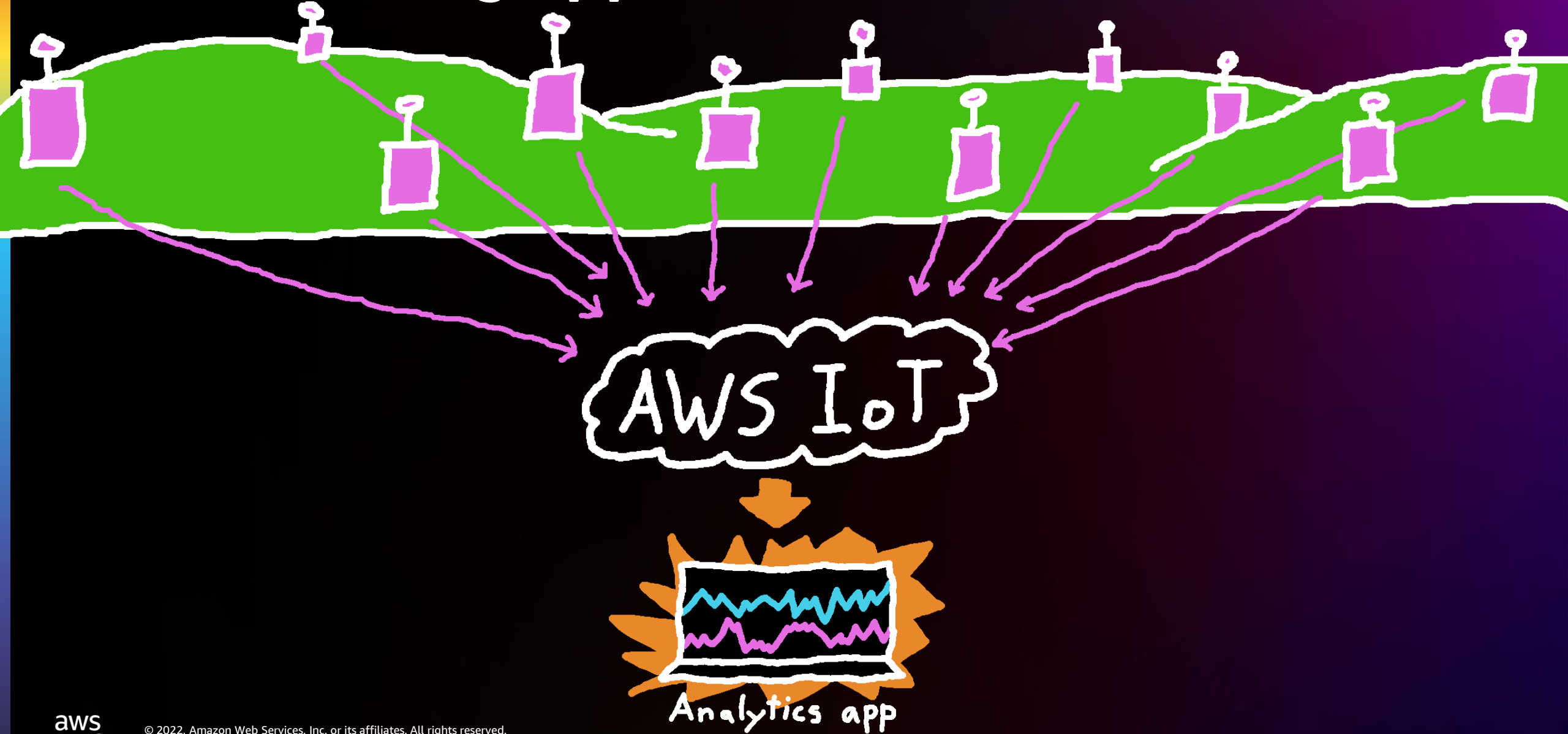
# Finding the root cause

- Propagate trace information
- Build lots of dashboards
- Get multi-dimensional
- Use high cardinality metrics
- Dig even deeper with log analysis
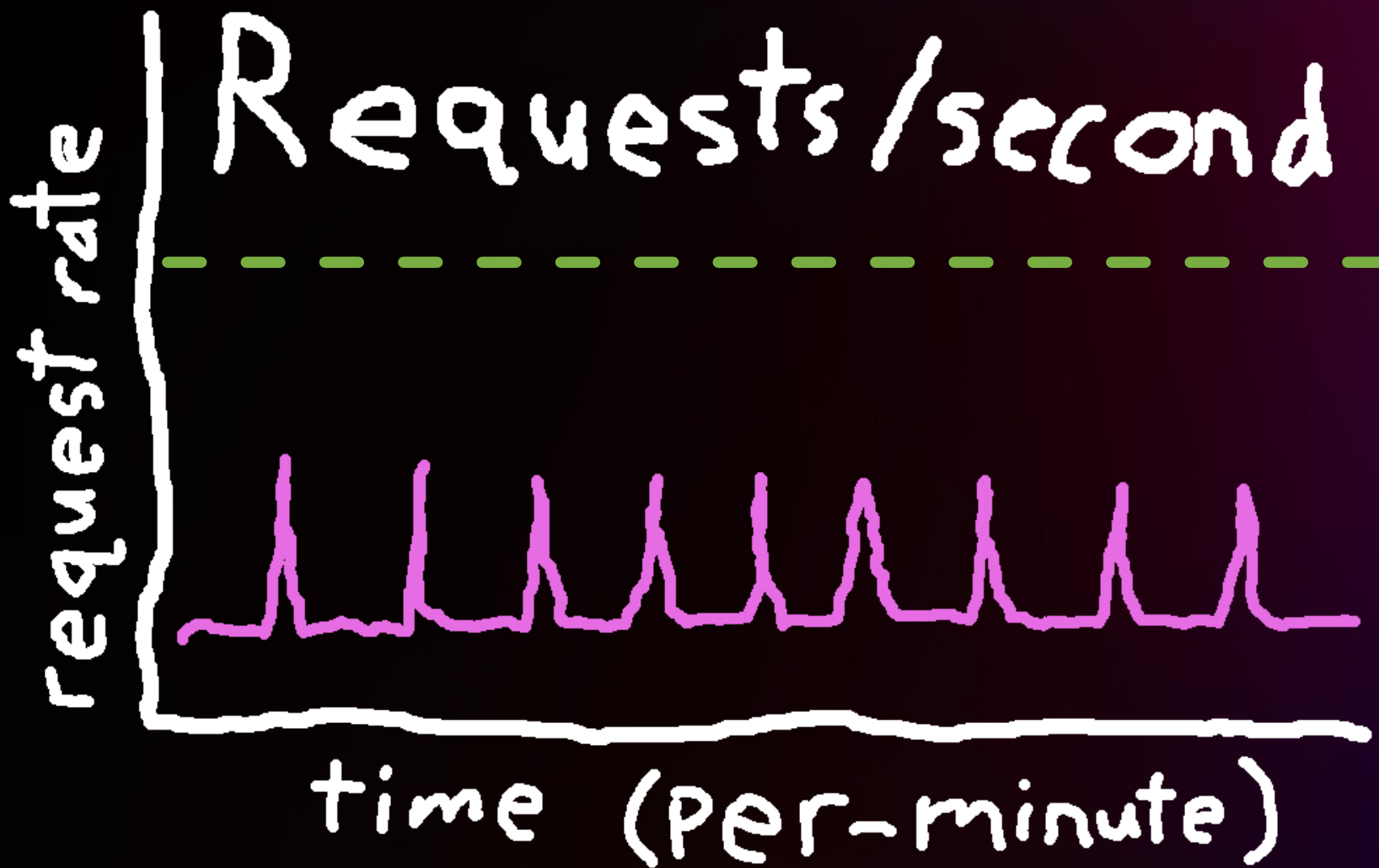- **Look at the raw logs**
- Analyze using profiler

ERROR

# Sifting through logs

# Sifting through logs

# Sifting through logs

# Raw logs

"I'm seeing errors talking to Amazon DynamoDB."

"Did the request make it to DynamoDB?"

```
java.net.UnknownHostException: dynamodb.us-east-1.amazonaws.com:
unknown error
   at java.net.Inet6AddressImpl.lookupAllHostAddr(Native Method)
   at java.net.InetAddress$2.lookupAllHostAddr(InetAddress.java:928)
   at java.net.InetAddress.getAddressesFromNameService(InetAddress.java:1323)
   at java.net.InetAddress.getAllByName0(InetAddress.java:1276)
   at java.net.InetAddress.getAllByName(InetAddress.java:1192)
   at java.net.InetAddress.getAllByName(InetAddress.java:1126)
   at com.amazonaws.SystemDefaultDnsResolver.resolve(SystemDefaultDnsResolver.j
```

# Raw logs

"I'm seeing errors talking to Amazon DynamoDB."

"What did DynamoDB say back?"

```
com.amazonaws.services.dynamodbv2.model.ConditionalCheckFailedException:
The conditional request failed (Service: AmazonDynamoDBv2; Status Code: 400;
Error Code: ConditionalCheckFailedException; Request ID: ...)
    at com.amazonaws.http.AmazonHttpClient.handleErrorResponse(AmazonHttpClient.
    at com.amazonaws.http.AmazonHttpClient.executeOneRequest(AmazonHttpClient.ja
    at com.amazonaws.http.AmazonHttpClient.executeHelper(AmazonHttpClient.java:4
    at com.amazonaws.http.AmazonHttpClient.execute(AmazonHttpClient.java:302)
    at com.amazonaws.services.dynamodbv2.AmazonDynamoDBClient.invoke(AmazonDynam
    at com.amazonaws.services.dynamodbv2.AmazonDynamoDBClient.updateItem(AmazonD
```

# Raw logs

"I'm seeing errors talking to Amazon DynamoDB."

"What happened right before that?"

```
[WARN]  2022-11-14T20:37:58.617Z    DynamoDB call timed, out; retrying 1/3
```

# Finding the root cause

- Propagate trace information
- Build lots of dashboards
- Get multi-dimensional
- Use high cardinality metrics
- Dig even deeper with log analysis
- Look at the raw logs
- **Analyze using profiler**

# Example: Inefficient code

```java
public ProductInfo getProductInfo(String customerId, String productId) {

    ProductInfo info = cache.get(customerId, productId);

    // compute pi just to make sure
    // TODO refactor or optimize
    for (int i = 0; i < 10000; i++) {
        calculatePI();
    }

    if (info == null) {
        info = db.query(customerId, productId);
    }

    return info;
}
```

No timer around this part of the code

# {CodeGuru} DemoProfilingGroup-WithIssues DEMO

| Data | View | Time range | | |
|---|---|---|---|---|
| CPU ▾ | Overview ▾ | 2021-07-16 @ 16:30 – 17:30 PDT | latest 12h ▾ | Actions ▾ |

Scroll up to see more data ▲

🔍 **4 Recommendations**

▼ **Legend**

My code ■
Other code ■

▼ **Minimap**



🔍 Reset zoom

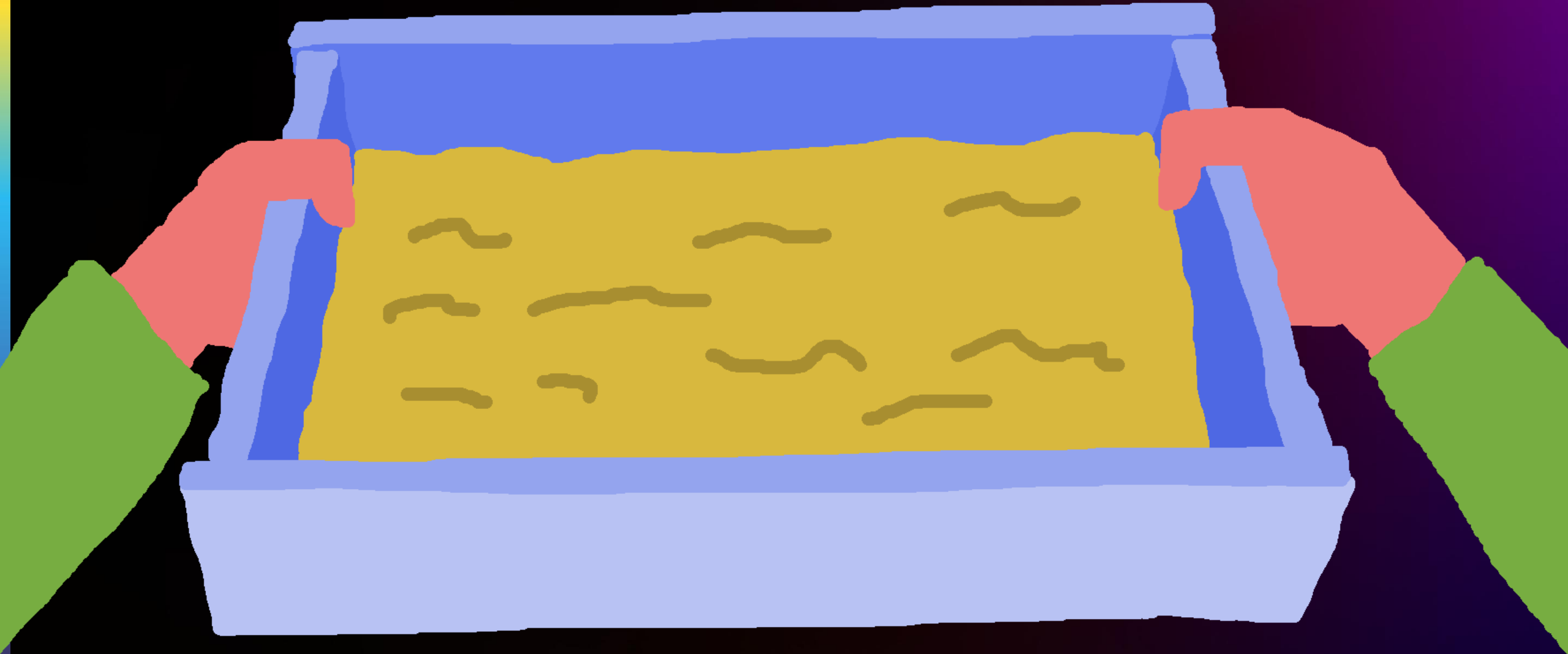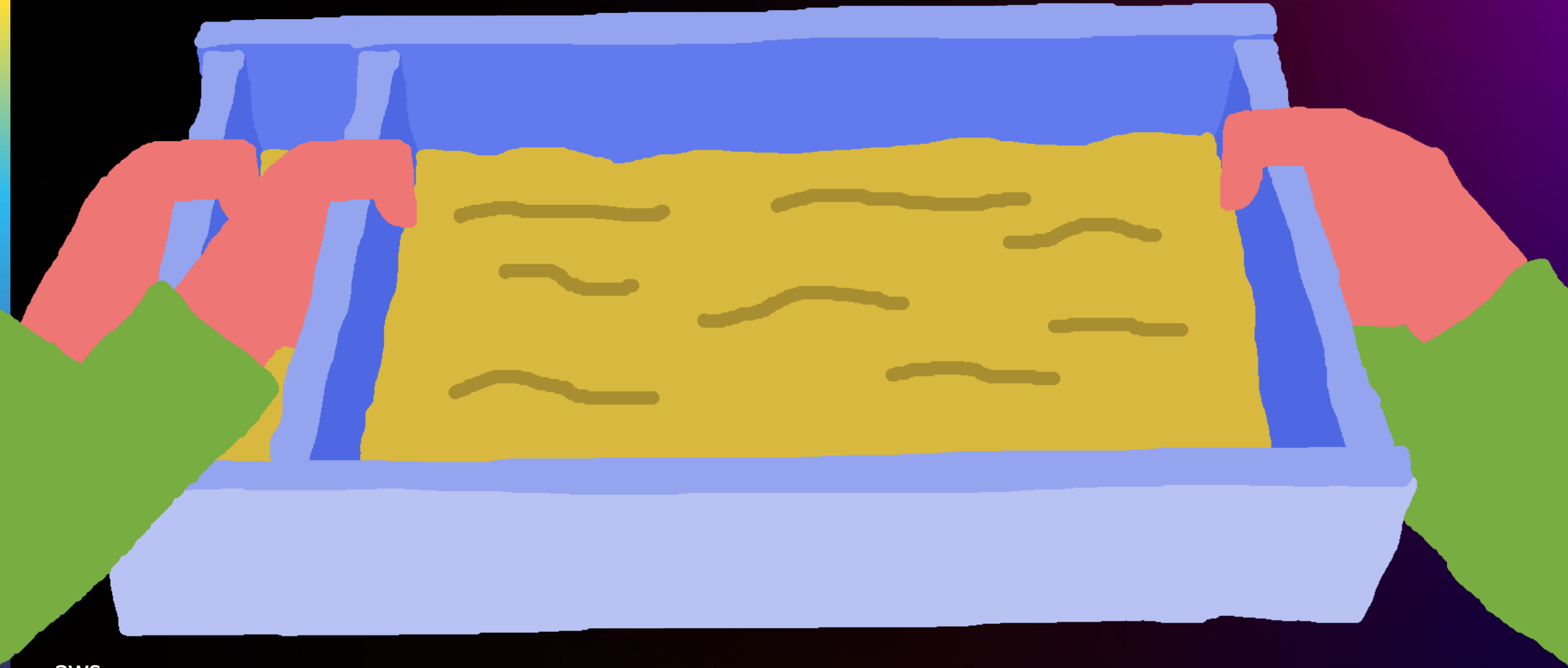| | | | | | | | | | | FileOutpu |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | FileOutpu |
| | | | | | | | | | | BufferedO |
| | | | | | | | | | | BufferedO |
| | | | | | | | | | | PrintStrea |
| | | | ZipFile.ge | | | | | Filter | FilterPrin | |
| | | | ZipFile.ge | | | | | Filter | FilterPrin | |
| | | | JarFile.ge | | | | | CloseS | CloseShiel | |
| | | | JarFile.ge | | | | | Output | OutputStre | |
| | | | URLClassPa | ZipFile.get | | | | Output | OutputStre | |
| | | | URLClassPa | ZipFile.get | ZipFile.getE | ZipFi | | Output | OutputStre | |
| | | | URLClassPat | JarFile.get | ZipFile.getEn | ZipFi | | Abstrac | AbstractOutpu | |
| | | | URLClassLoa | JarFile.get | JarFile.getEn | JarFi | | Abstrac | AbstractOutpu | |
| | MainCli | | URLClassLoa | URLClassPat | JarFile.getJa | JarFi | | Abstrac | AbstractOutpu | |
| | Protocol | | AccessContro | URLClassPat | URLClassPath$ | URLCl | | Appende | AppenderContr | |
| | Internal | | URLClassLoad | URLClassPat | URLClassPath$ | URLCl | Annot | Appende | AppenderContr | |
| | Closeabl | | ClassLoader.g | URLClassPat | URLClassPath. | URLCla | Annot | Appende | AppenderContr | |
| | Closeabl | | URLClassLoader | URLClassLoa | URLClassLoade | URLCla | Annot | Appende | AppenderContr | |
| Amazon | SdkHttpC | | VersionInfo.loadVer | URLClassLoa | URLClassLoade | URLCla | Annot | LoggerC | LoggerConfig. | |
| AmazonHttpClient | ApacheCo | | VersionInfo.getUser | AccessContr | AccessControl | Access | POJOProp | POJOPro | LoggerC | LoggerConfig. | |
| AmazonHttpClient$ | ApacheCo | HttpClientBuilder.build | | URLClassLoa | URLClassLoade | URLCla | POJOPropertiesCol | LoggerC | LoggerConfig. | |
| AmazonHttpClient$ | ApacheHttpClientFactory.create | | URLClassLoa | ClassLoader.g | ClassLoad | POJOPropertiesCol | LoggerC | LoggerConfig. | Deflater.deflateByt |
| AmazonHttpClient$ | ApacheHttpClientFactory.create | | CompoundEnum | URLClassLoade | URLClassLoa | BasicBeanDescript | LoggerC | LoggerConfig. | Deflater.deflate |
| AmazonHttpClient$ | AmazonHttpClient.<init> | | CompoundEnum | Class.getReso | Class.getRe | BasicSerializerFa | LoggerC | LoggerConfig. | Deflater.deflate |
| AmazonHttpClient$ | AmazonHttpClient.<init> | | Collections. | HandlerChainF | HandlerChainF | BeanSerializerFacto | LoggerC | LoggerConfig. | IDATOutputStream.de |
| AmazonHttpClient. | AmazonWebServiceClient.<init> | | HandlerChain | HandlerChainF | HandlerChainF | BeanSerializerFacto | AwaitCo | AwaitCompleti | IDATOutputStream.wri |
| AmazonHttpClient. | AmazonWebServiceClient.<init> | AmazonSQSClient.init | BeanP | SerializerProvider. | Logger. | Logger.log | PNGImageWriter.encod |
| AmazonSQSClient.d | AmazonSQSClient.<init> | | BeanS | SerializerProvider. | Abstrac | AbstractLogge | PNGImageWriter.write |
| AmazonSQSClient.i | AmazonSQSClient.<init> | | BeanS | SerializerProvider. | AbstractLo | AbstractLogger | PNGImageWriter.write |
| AmazonSQSClient.i | AmazonSQSClientBuilder.build | | Defau | SerializerProvider. | AbstractLo | AbstractLogger | ImageWriter.write |
| AmazonSQSClient.e | AmazonSQSClientBuilder.build | | DefaultSerializerProvider | AbstractLo | AbstractLogger | ImageIO.doWrite |
| AmazonSQSClient.e | AwsSyncClientBuilder.build | | ObjectMapper._configAndWr | AbstractLo | AbstractLogger | ImageIO.write | TaskP |
| AmazonSQSClient.r | Main.sqsClient | | ObjectMapper.writeValueAs | AbstractLo | AbstractLogger | ImageProcessor$Brighten | Main. |
| | ImageProcessor.extractTasks | | | | | | ImageProcessor$Brighten | Main$ |
| ImageProcessor.run | | | | | | | | Execu |
| Main$Lambda.run | | | | | | | | Futur |
| Executors$RunnableAdapter.call | | | | | | | | Sched |
| FutureTask.run | | | | | | | | Sched |
| ThreadPoolExecutor.runWorker | | | | | | | | |
| PS_Scaveng | ThreadPoolExecutor$Worker.run | | | | | | | |
| GarbageCol | Final | Thread.run | | | | | | ProfilingComman |
| ALL | | | | | | | | |

# Finding the root cause

- Propagate trace information
- Build lots of dashboards
- Get multi-dimensional
- Use high cardinality metrics
- Dig even deeper with log analysis
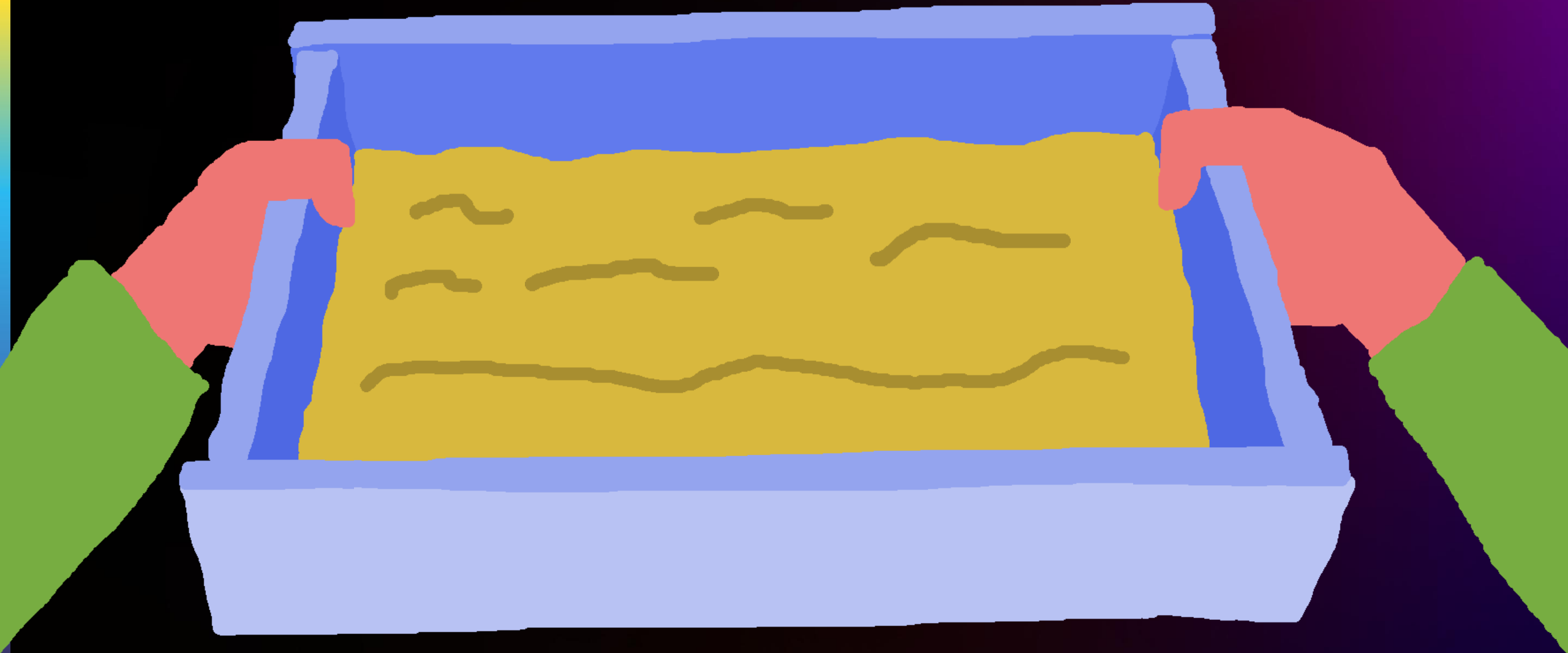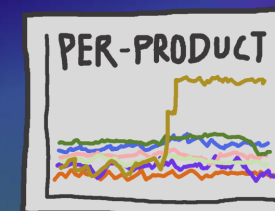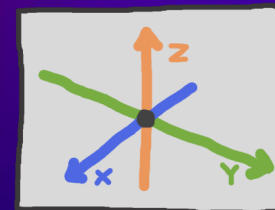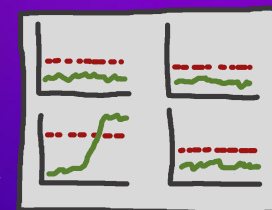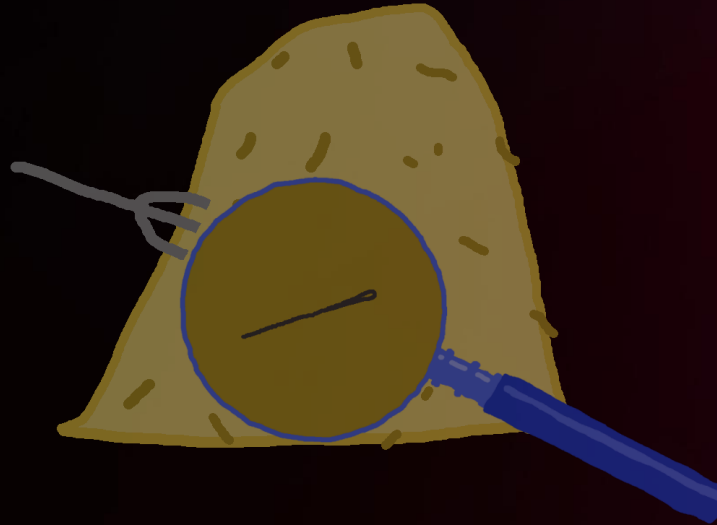- Look at the raw logs
- Analyze using profiler

PER-PRODUCT

SELECT
FROM
WHERE

ERROR

# Chapter three

The DevOps
flywheel at Amazon

Find the
root cause

Measure from
everywhere

# Measure from everywhere

- Be your own customer
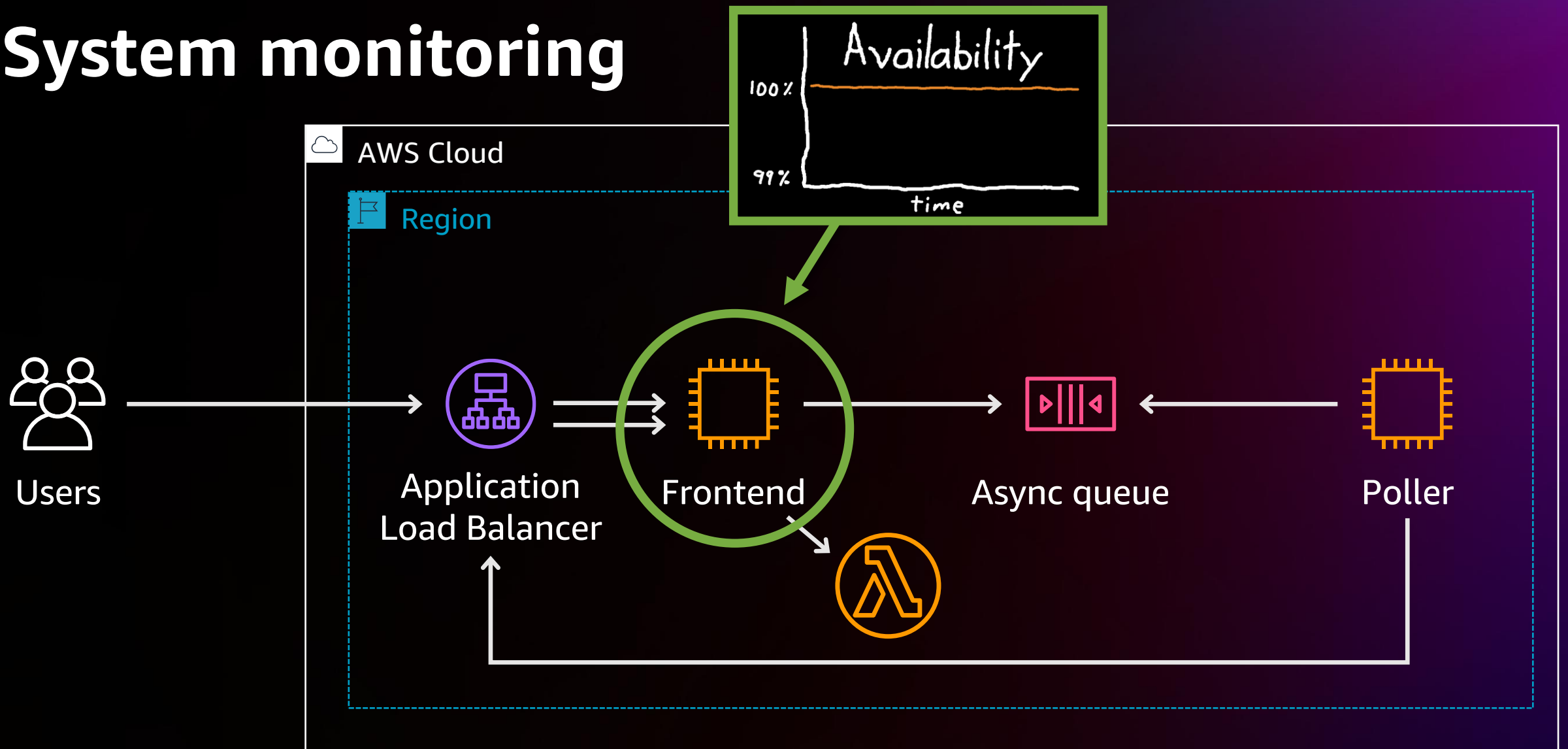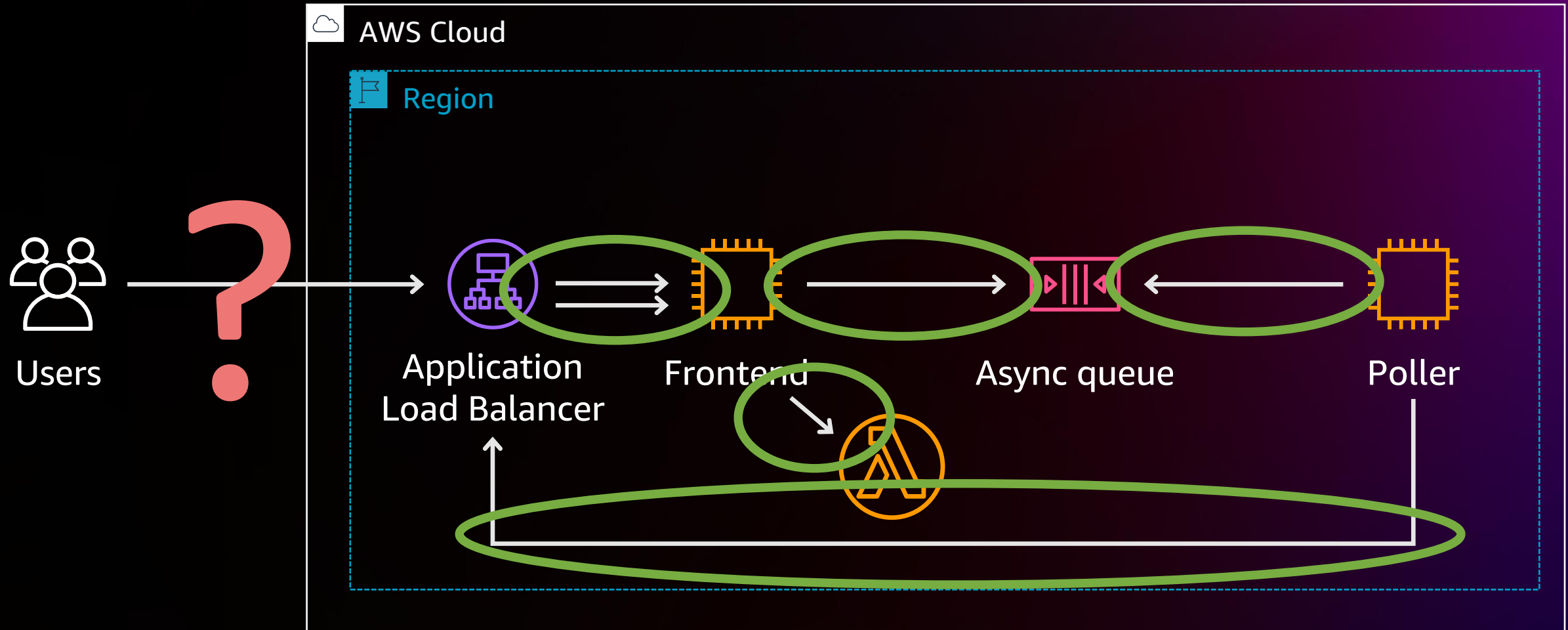- Instrument close to the customer
- Alarm on everything

# Measure from everywhere

- **Be your own customer**
- Instrument close to the customer
- Alarm on everything

# System monitoring



Availability

100%

99%

time

AWS Cloud

Region

Users

Application
Load Balancer

Frontend

Async queue

Poller

# System monitoring

# Solving browser testing

Diverse device types, browsers

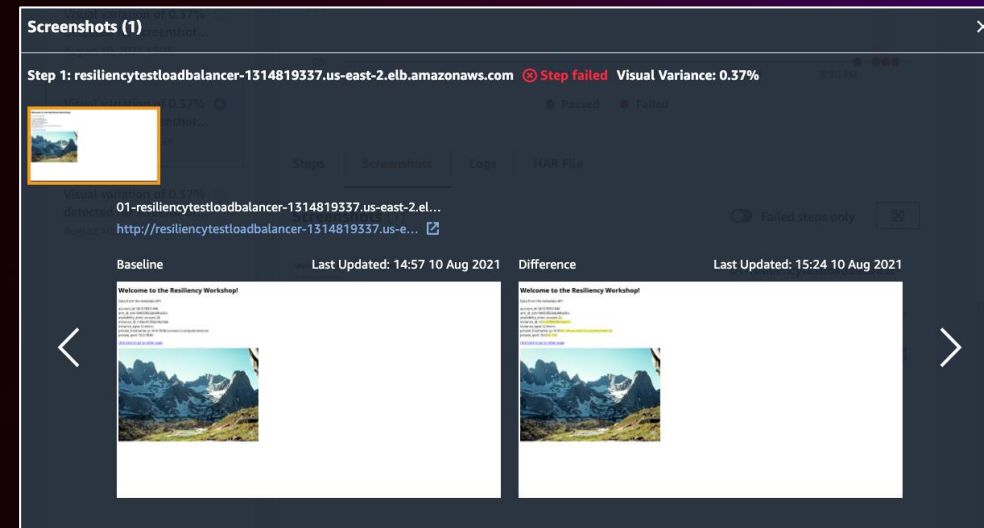Testing flows of
human interaction

# CloudWatch Synthetics

Pre-built blueprint scripts
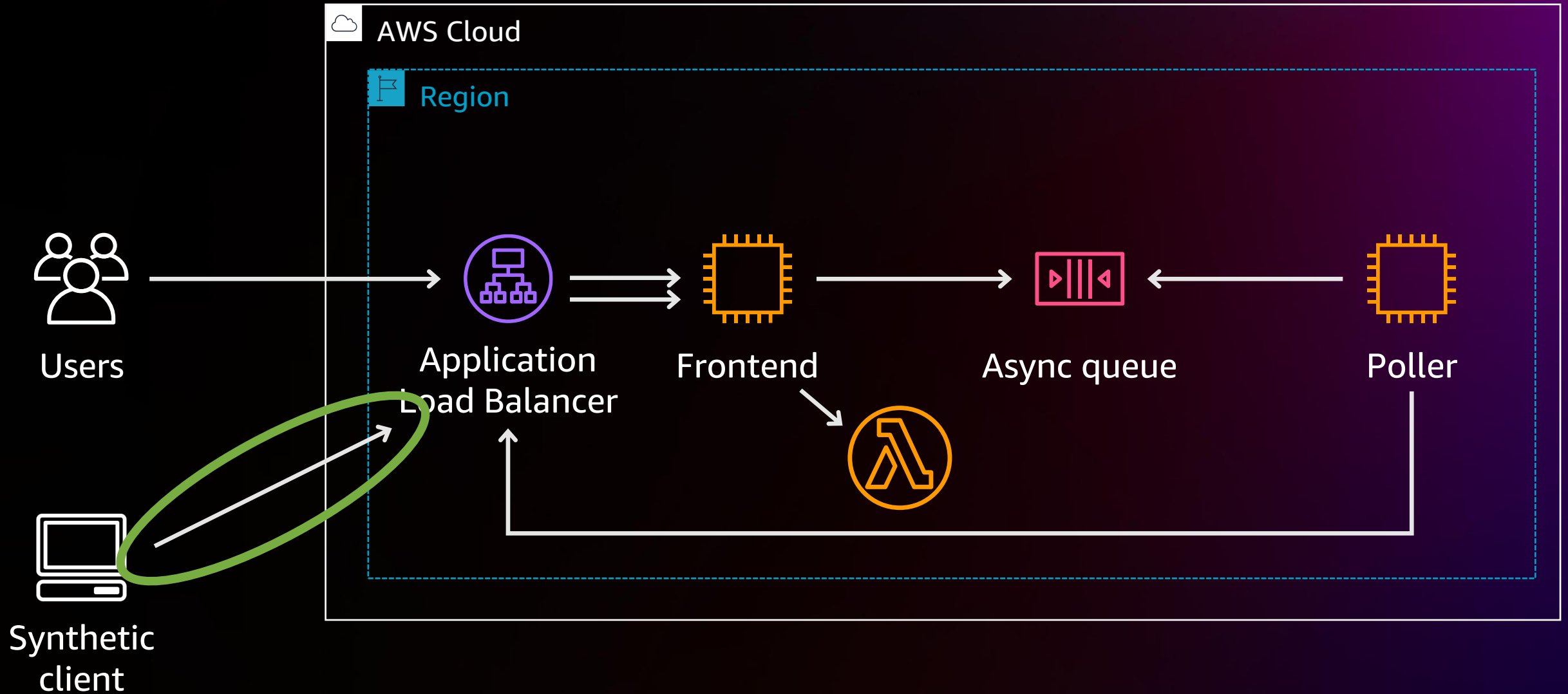& visual script reorder
(chrome plugin)

Visually inspect
synthetic failures

Fully customizable scripting in
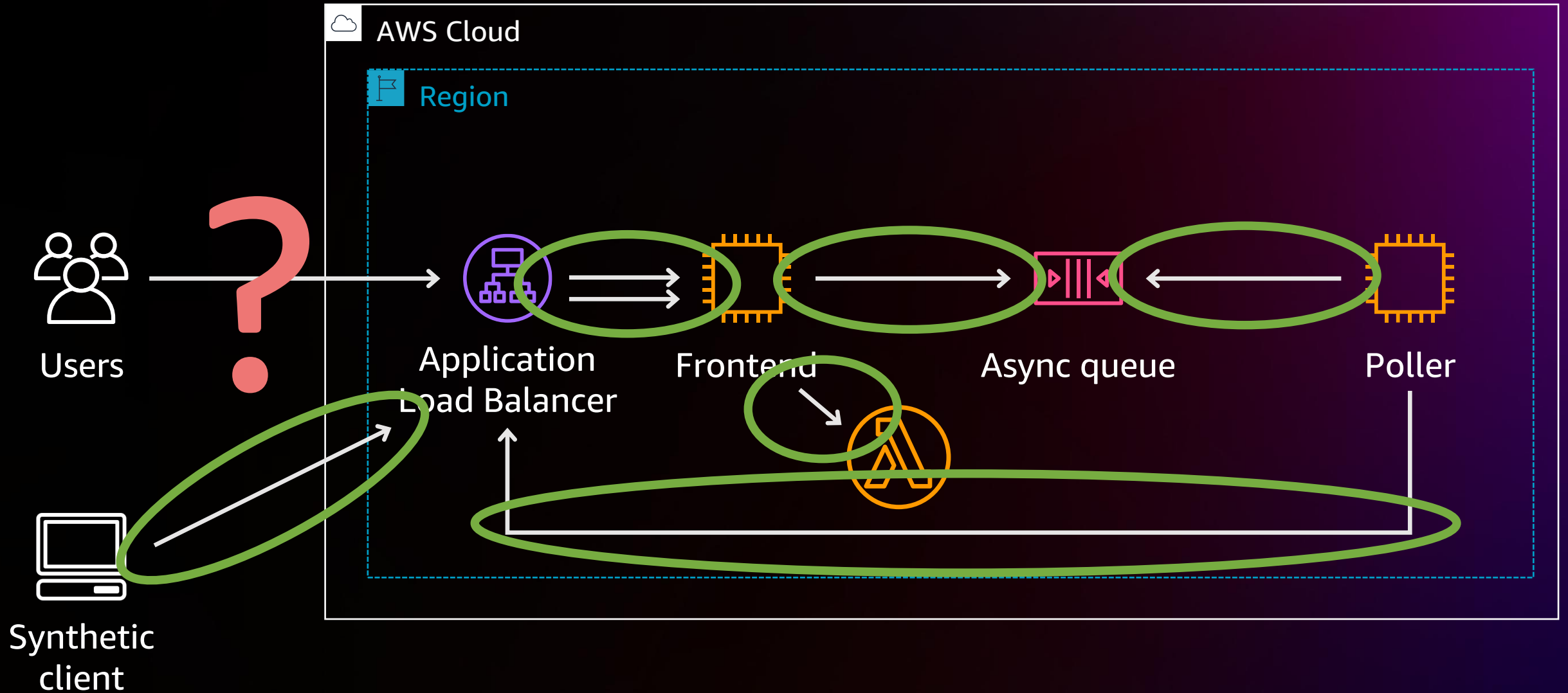NodeJS and Python using
Selenium & Puppeteer

# System monitoring

# Measure from everywhere

- Be your own customer
- **Instrument close to the customer**
- Alarm on everything

# System monitoring



Users

Synthetic
client

AWS Cloud

Region

Application
Load Balancer

Frontend
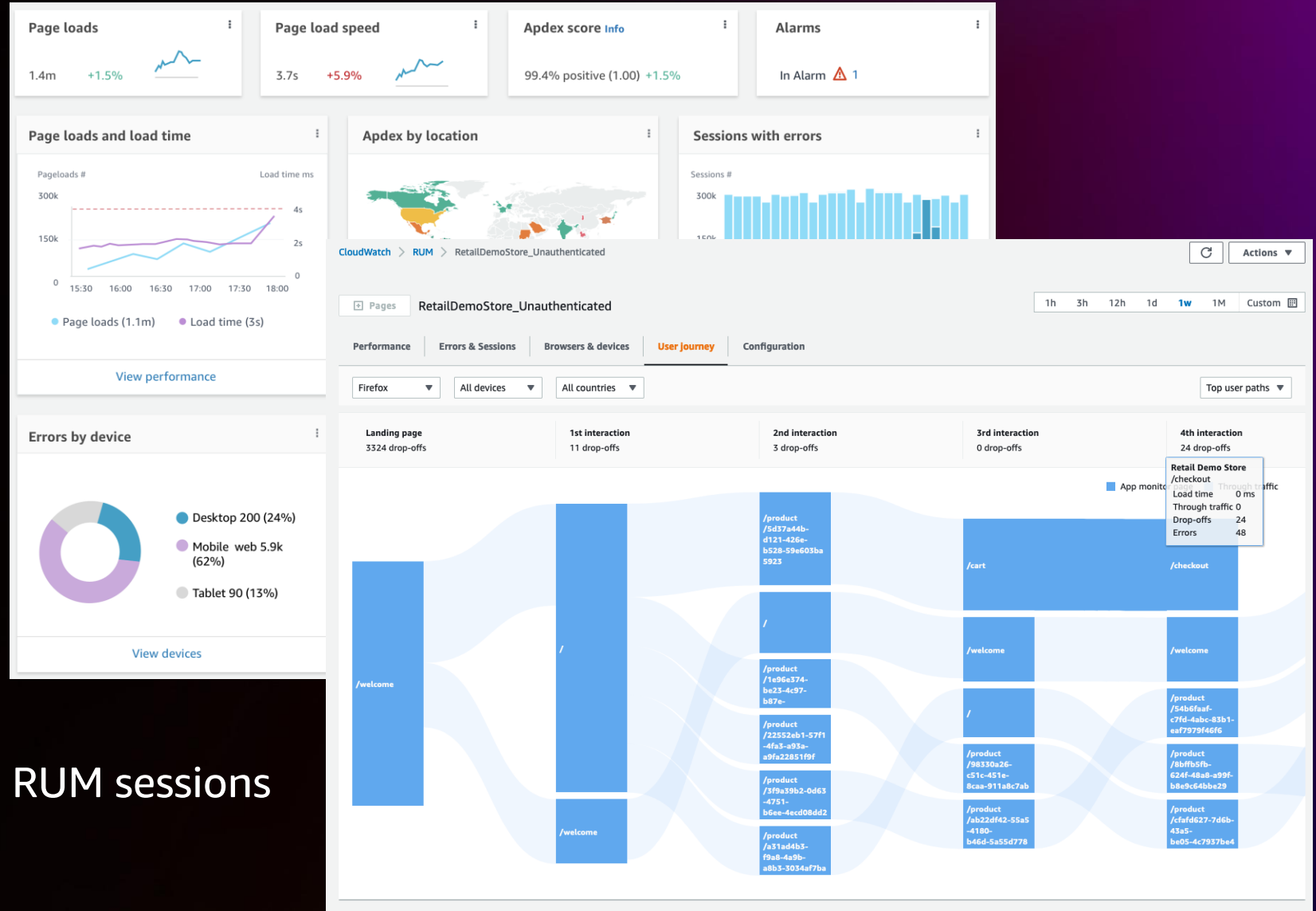
Async queue

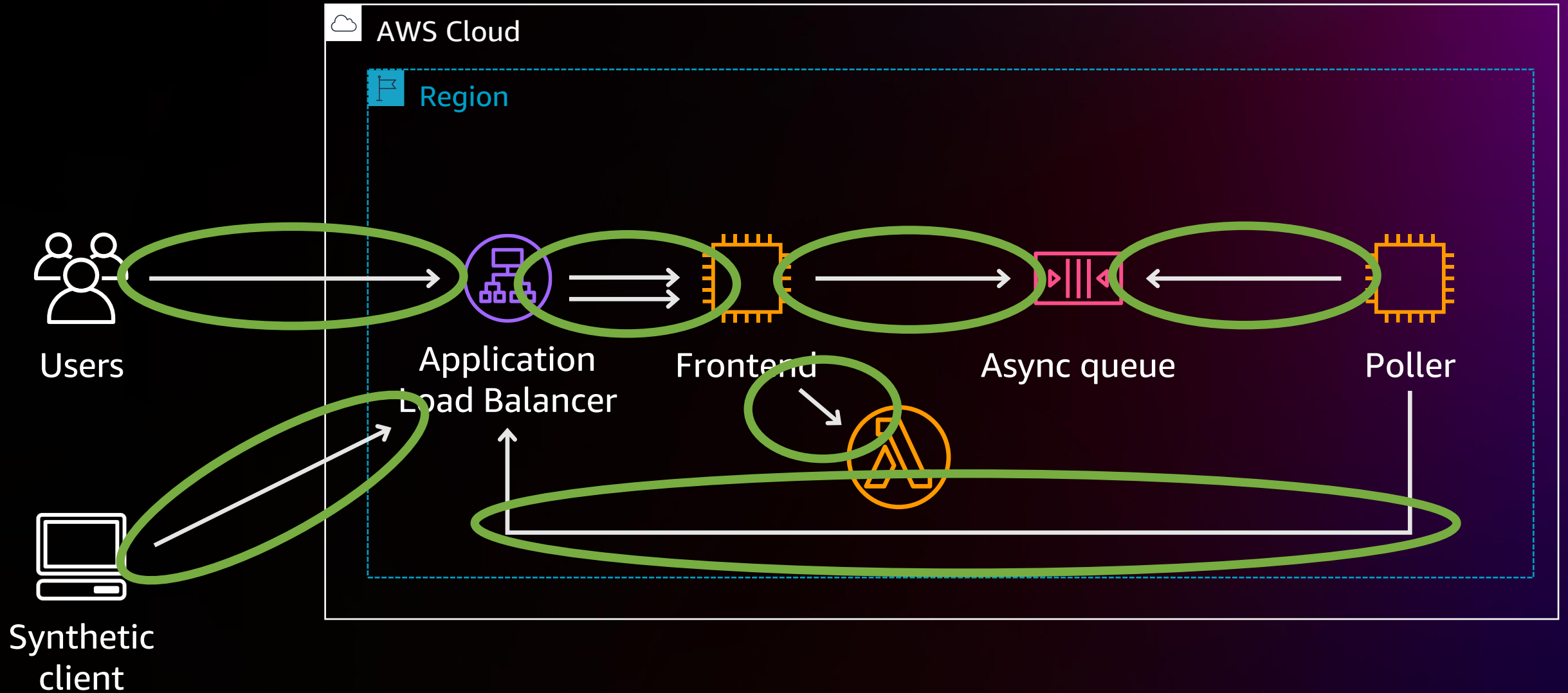Poller

# CloudWatch RUM

**Quantifying slowness**

- Apdex scores
- Errors by sessions, devices, canaries, regions, pages
- Core web vitals
- Largest contentful paint
- First input delay
- Cumulative layout shift

**Debugging slowness**

- Page load steps
- Resource requests
- View X-Ray traces for errors in RUM sessions
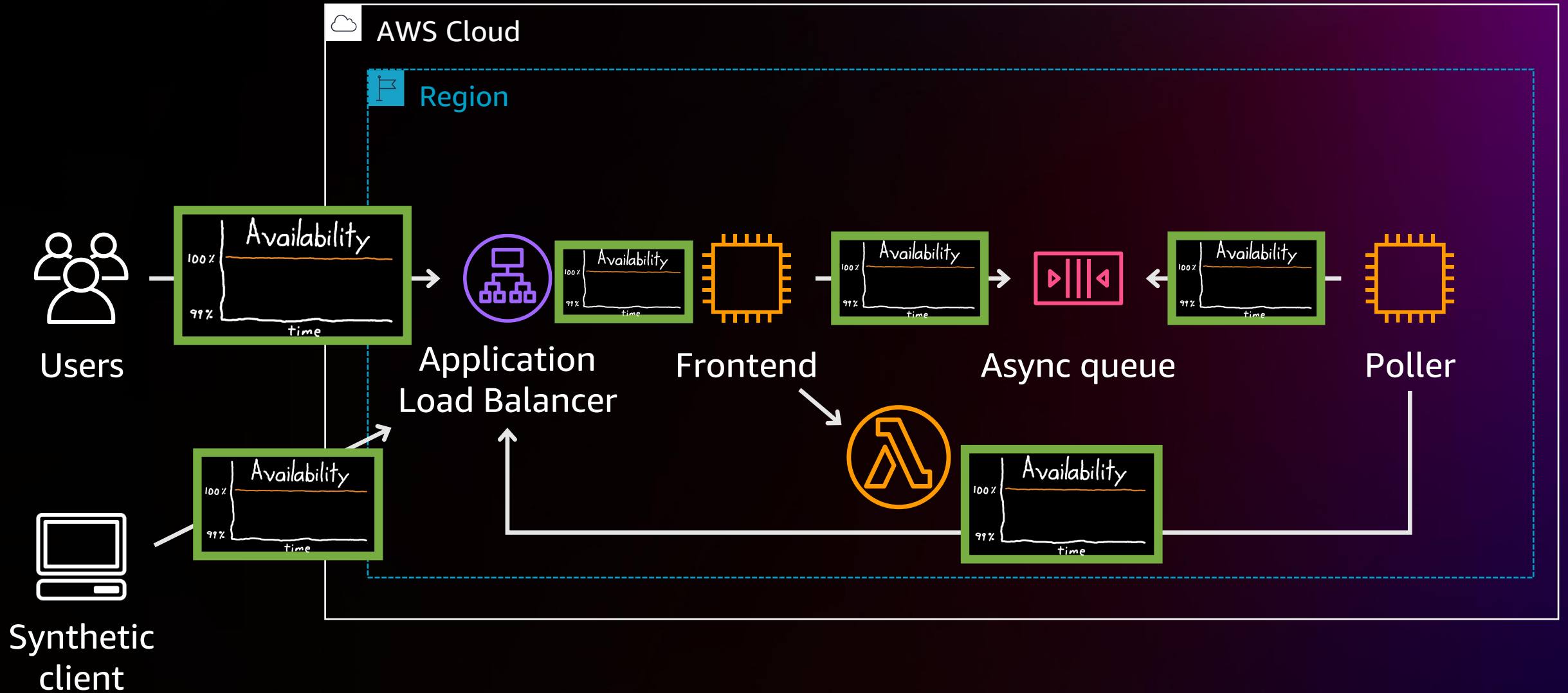- JavaScript, HTTP errors

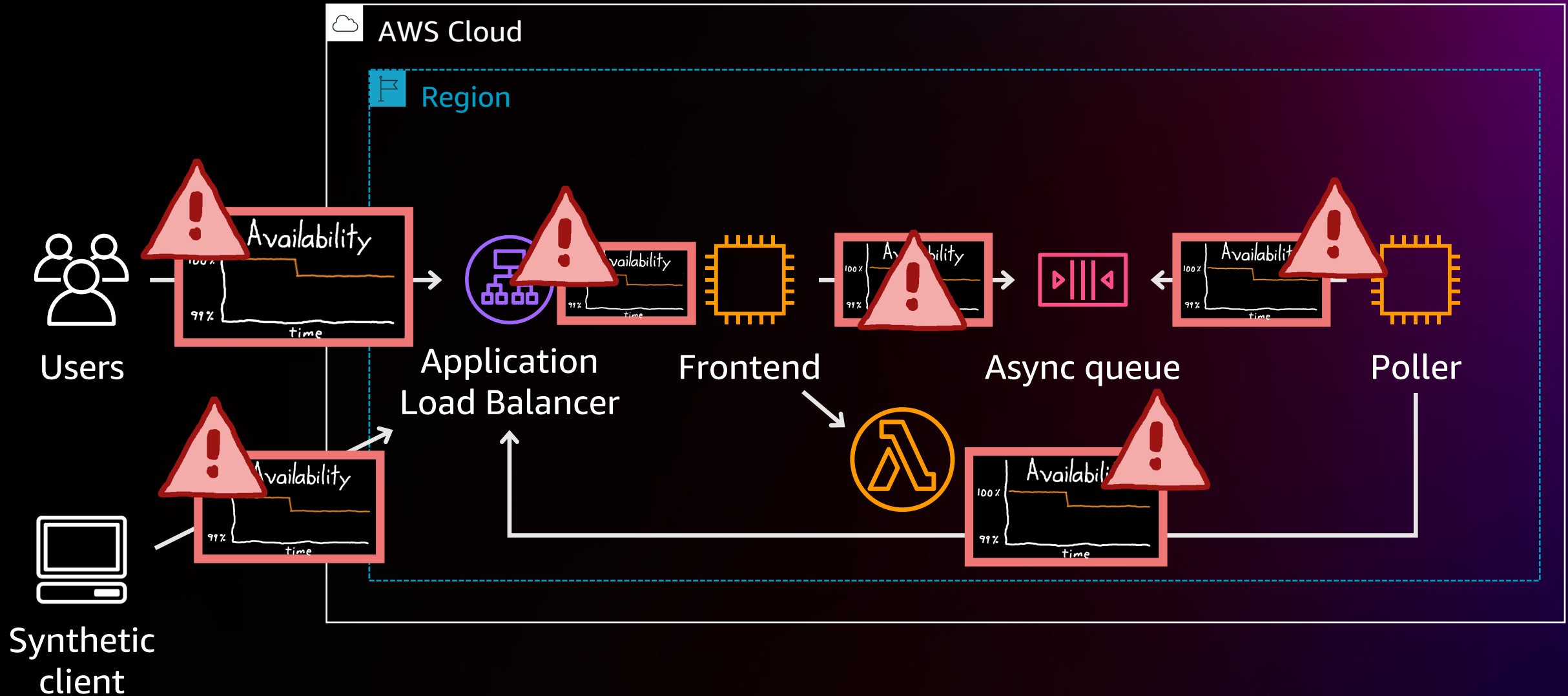# System monitoring

# Measure from everywhere

- Be your own customer
- Instrument close to the customer
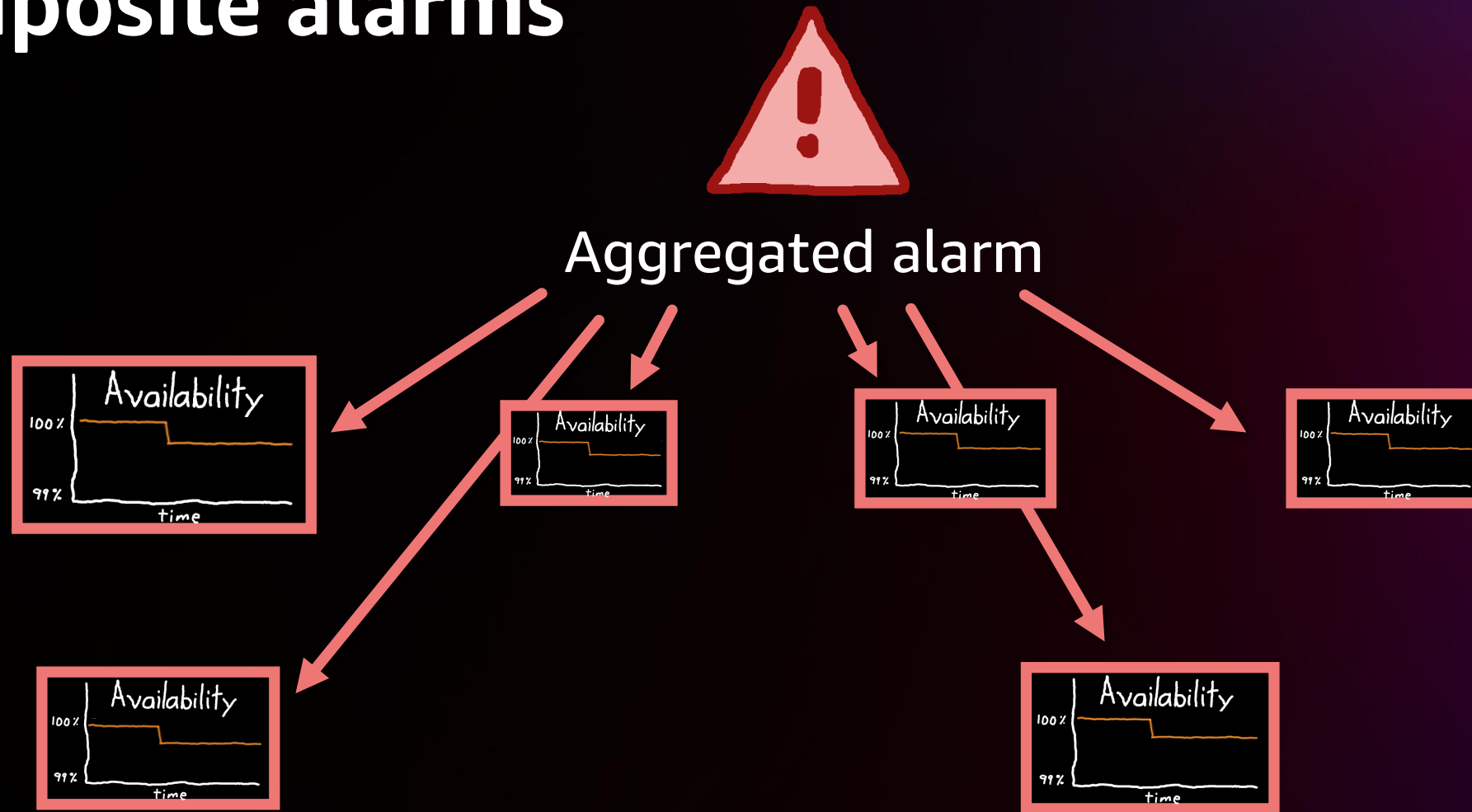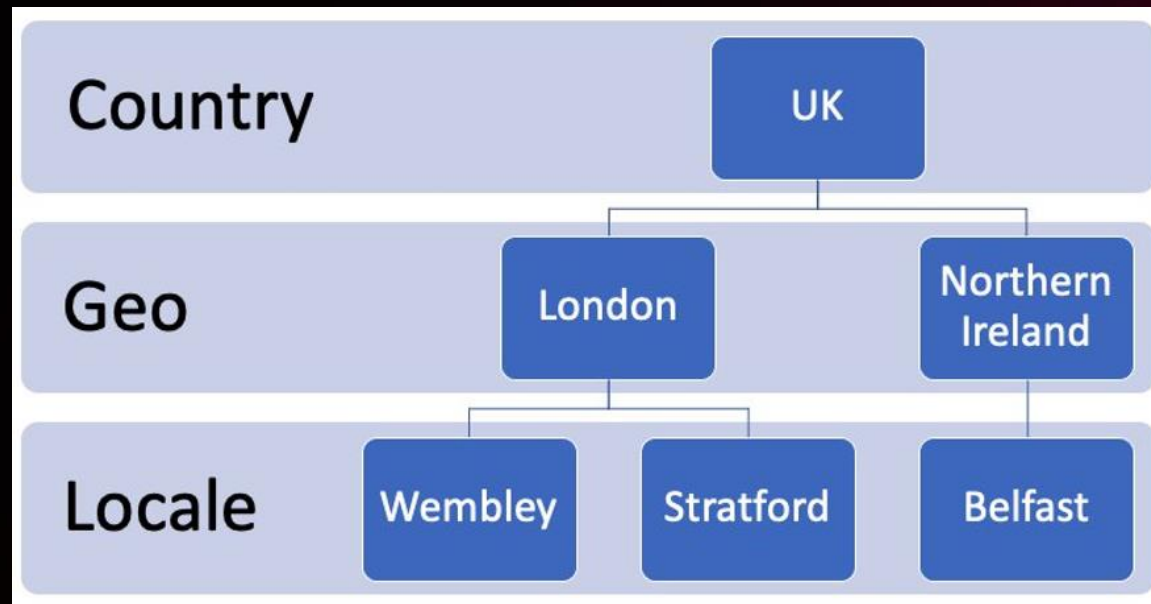- **Alarm on everything**

# System monitoring



© 2022, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# System monitoring



Users

Synthetic client

AWS Cloud

Region

Application Load Balancer

Frontend

Async queue

Poller

# Composite alarms



Aggregated alarm

# Reducing noise at scale – BT



https://aws.amazon.com/blogs/mt/how-bt-uses-amazon-cloudwatch-to-monitor-millions-of-devices/
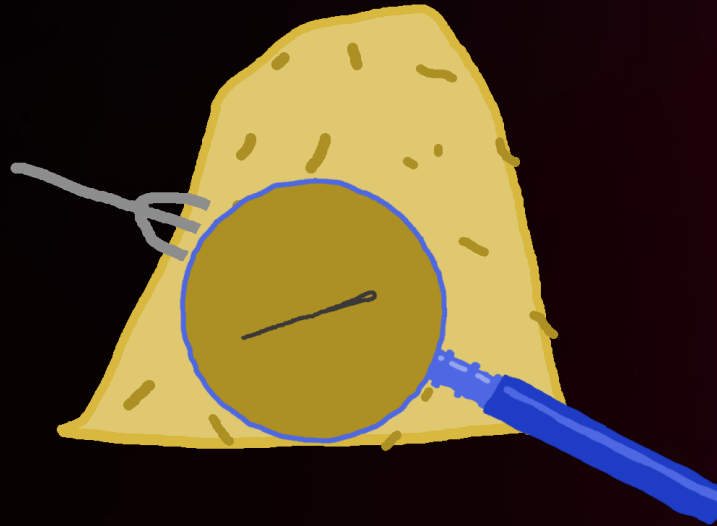
# Measure from everywhere

- Be your own customer
- Instrument close to the customer
- Alarm on everything
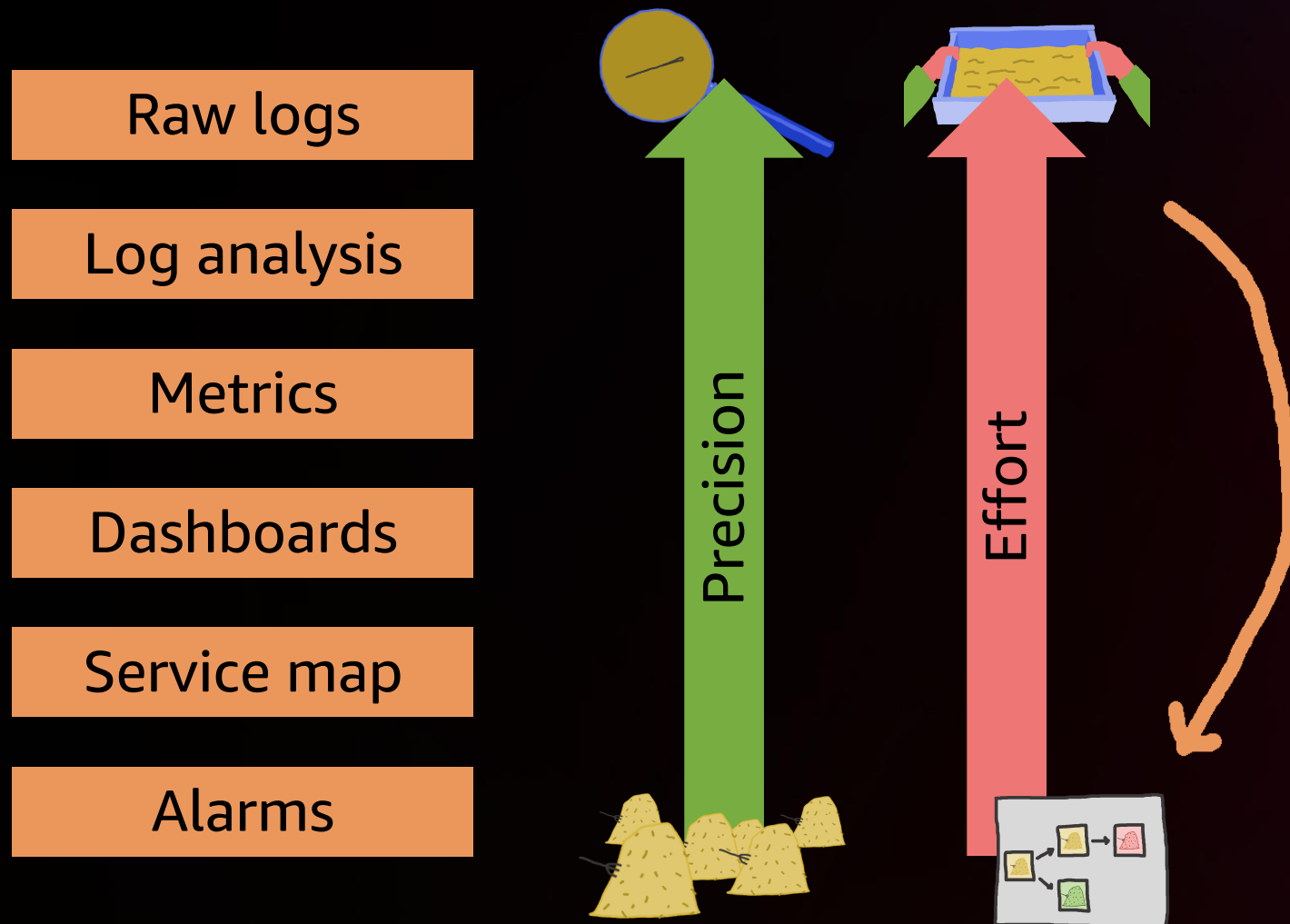
# Key takeaways

The DevOps
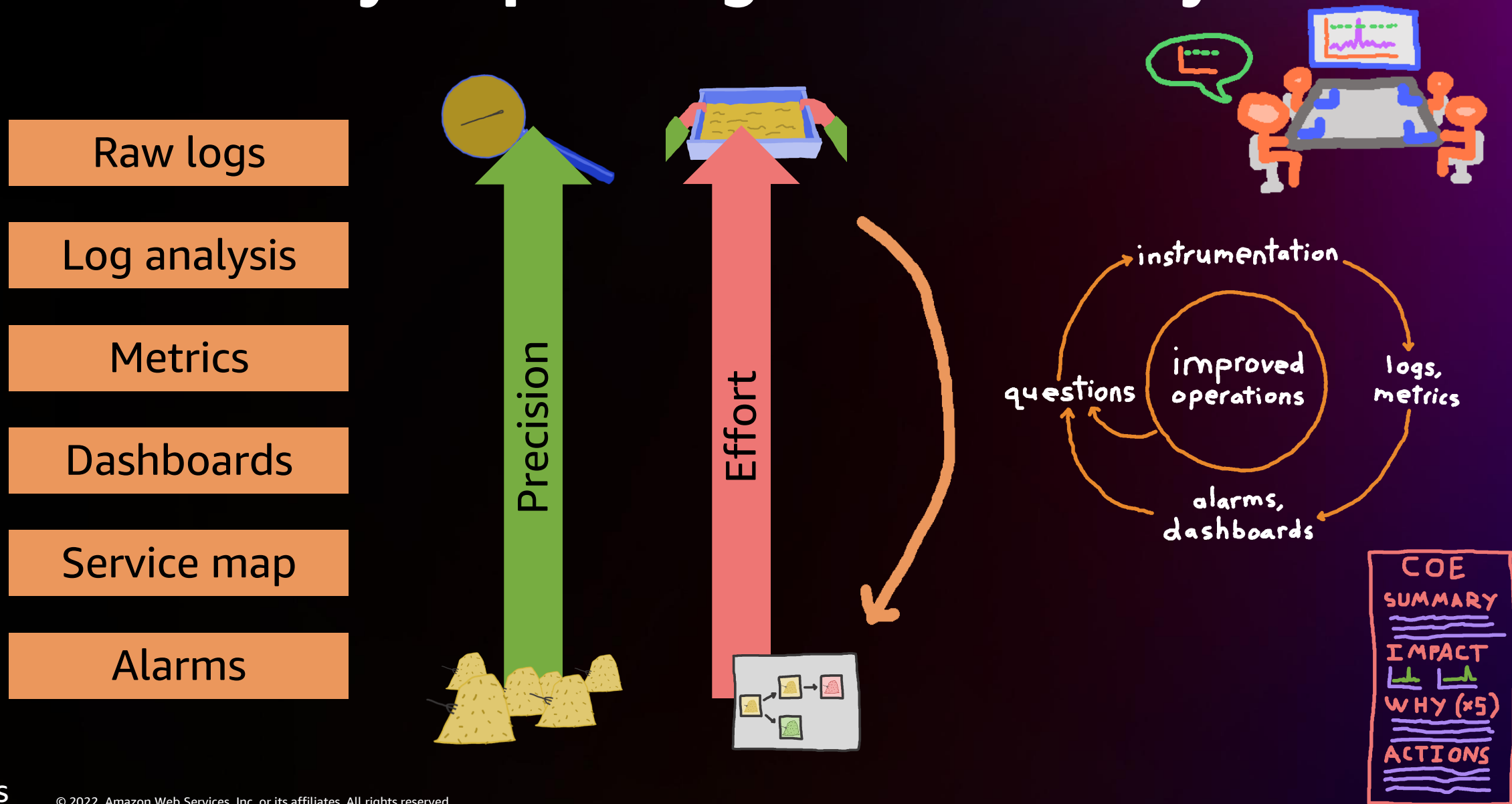flywheel at Amazon

Find the
root cause

Measure from
everywhere

# Continuously improving observability

Raw logs

Log analysis

Metrics

Dashboards

Service map

Alarms

Precision

Effort

# Continuously improving observability



Raw logs

Log analysis

Metrics

Dashboards

Service map

Alarms

Precision

Effort

instrumentation

questions

improved operations

logs, metrics

alarms, dashboards

COE SUMMARY

IMPACT

WHY (×5)

ACTIONS

# Thank you!

David Yanacek

Ian McGarry

@dyanacek

ianmcgarry

https://aws.amazon.com/builders-library/

Please complete the session survey in the **mobile app**