

Introducing The Amazon Builders' Library

Andrew Certain

Sr. Principal Engineer
Database Services

Becky Weiss

Sr. Principal Engineer
AWS Identity

Colm MacCárthaigh

Sr. Principal Engineer
EC2 Networking

David Yanacek

Principal Engineer
AWS Lambda

Marc Brooker

Sr. Principal Engineer
AWS Serverless Applications

AWS re:Invent

© 2019, Amazon Web Services, Inc. or its affiliates. All rights reserved.



The Amazon Builders' Library is a collection of living articles that are written by Amazon senior technical leaders and engineers and that describe how Amazon develops, architects, releases, and operates technology. In this session, five principal engineers share hard-learned lessons from their experiences building reliable services at Amazon. David Yanacek, Andrew Certain, Becky Weiss, Colm MacCarthaigh, and Marc Brooker each share a personal story highlighting a current best practice. The engineers discuss how Amazon uses timeouts, how we think about back-offs and retries, our approach to taking dependencies, how we measure performance, and how we use shuffle sharding.



Amazon
DynamoDB



AWS Lambda



Amazon Virtual
Private Cloud
(Amazon VPC)



Amazon Elastic
Block Store
(Amazon EBS)



Amazon
Route 53



Amazon
CloudFront



Elastic Load
Balancing



Amazon Quantum
Ledger Database
(Amazon QLDB)

Lightning talks

Monitoring with percentiles

Static stability

Shuffle sharding

Retries, backoff, and jitter

Meta lesson: Learn from each other

Lightning talks

A. Andrew

B. Becky

C. Colm

D. David

E. Marc

Monitoring with percentiles

<http://blog.tacertain.com/p-four-nines/>



"P-four-nines"

=

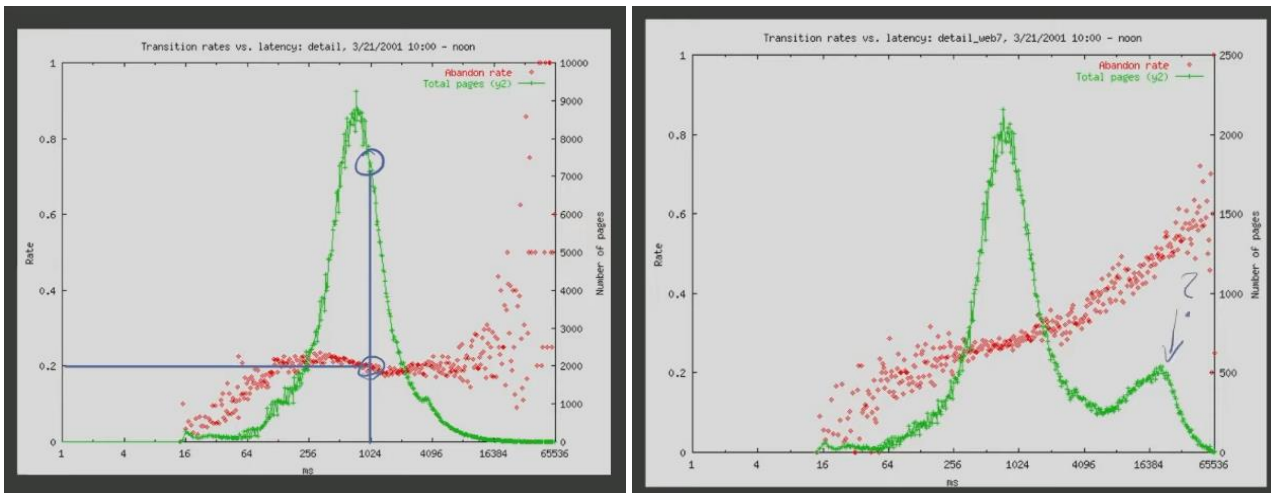
99.99th percentile

=

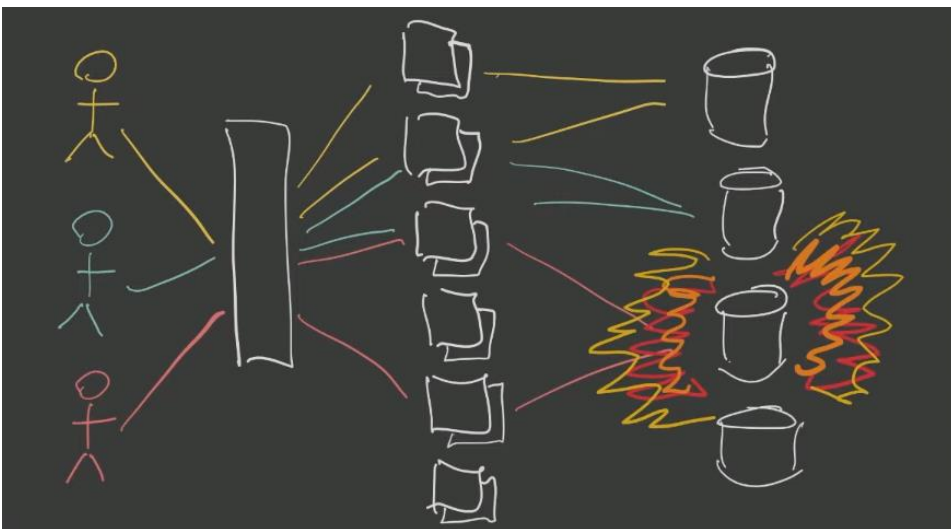
Value that is $\geq 99.99\%$ of
the data points



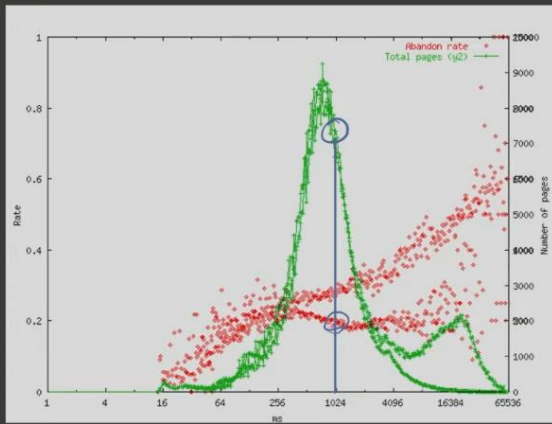
This is the 2001 version of the Amazon website. We switched from looking at the average values to using percentiles



About 7000 pages had 1 second latency. Less latency means that the abandon rate increases



This is a rendering of the website in 2001, LBs, web servers and the databases. One DB was having troubles and adding latencies to the page load times randomly for the page views.



The bad DB was causing the high abandon rates when the customers were viewing pages that interact with it. We determined that the user threshold is about 6 seconds for the page to load, we need to switch to percentiles and use the 99.99 percentiles.

Creating a Embedded Metric Format log in Lambda

Screenshot of Lambda function using the Node.js client library

```

1 const { metricScope } = require("aws-embedded-metrics");
2
3 const myFunction = metricScope(metrics =>
4   async (event, context) => {
5     const now = new Date().getTime();
6     metrics.putMetric("processingLatency", now - event.eventTime, "Milliseconds");
7     metrics.putDimensions({ functionVersion: process.env.AWS_LAMBDA_FUNCTION_VERSION });
8     metrics.setProperty("RequestId", context.awsRequestId);
9     metrics.setProperty("Device", context.sourceDevice);
10  });
11
12 exports.handler = myFunction;
13

```

<https://github.com/aws-labs/aws-embedded-metrics-node>

You can get percentiles data using a Lambda out of CloudWatch

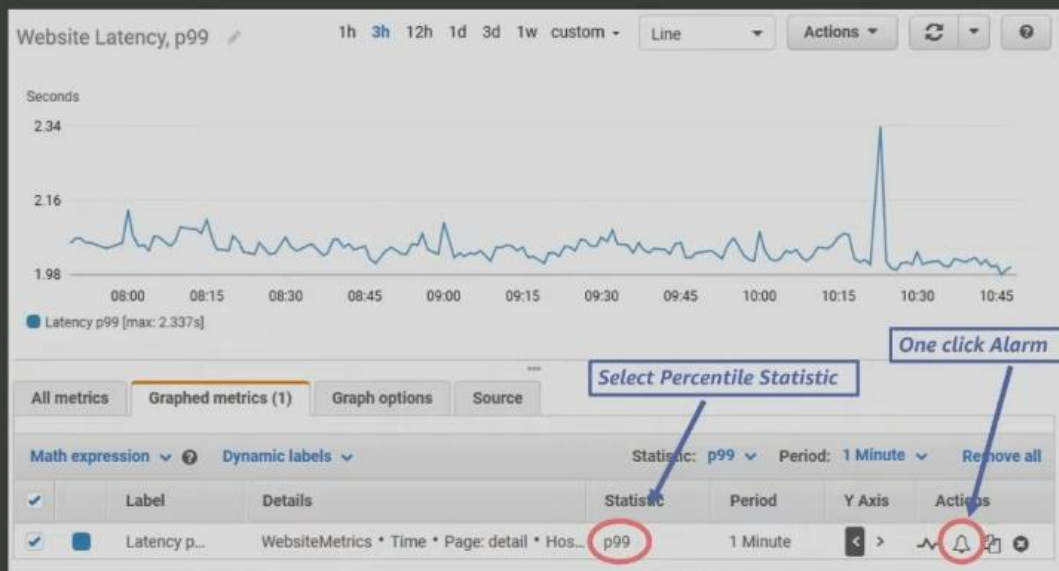
Example JSON log event

```
Time (UTC +00:00)  Message
2019-11-22
23:26:47  START RequestId: eeb90ca1-8ac1-40d5-a56d-84204ff23ec3 Version: $LATEST
23:26:47  2019-11-22T23:26:47.602Z eeb90ca1-8ac1-40d5-a56d-84204ff23ec3 INFO ["LogGroup":"CustomMetricsTest","ServiceName":"CustomMetricsTest","ServiceType":"AWS::Lambda::Function","FunctionVersion":"$LATEST","RequestID":"42201529-16f6-4a83-b8f8-fa3f99100f8","Device":"61279781-c6ac-46f1-baf7-22c88a78162","ExecutionEnvironment":"AWS::Lambda::Nodejs10.x","MemorySize":"1216","LogStreamId":"2019/11/22/[SLATEST]3bbddad832a047d2a6b00055a27cbb8","aws":{"Timestamp":"1574465267600","CloudWatchMetrics":{"Dimensions":["LogGroup","ServiceName","ServiceType","FunctionVersion"],"Metrics":[{"Name":"ProcessingLatency","Unit":"Milliseconds","Namespace":"aws-embedded-metrics"}]}}]
2019-11-22T23:26:47.682Z eeb90ca1-8ac1-40d5-a56d-84204ff23ec3 INFO {"ProcessingLatency": 229.68621283857877}
```

Environment and resource information included by default

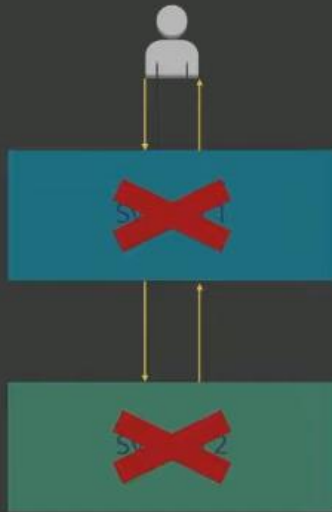
Embedding custom metrics in detailed log event – specify namespace, metric names, dimensions, and units

Viewing percentile information on a Amazon CloudWatch metric



**Static stability in AWS:
Resilience in the face of impaired
dependencies**

Dependencies: We all have them



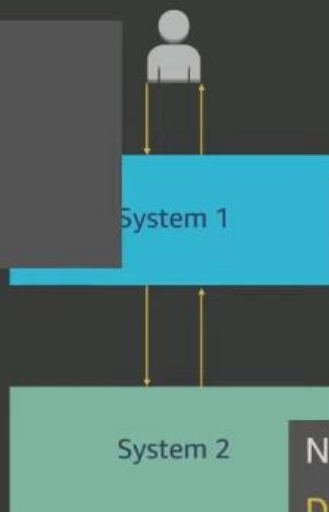
Dependencies: We all have them

Math:

For all $0 < p, q < 1$:

$$p * q < p$$

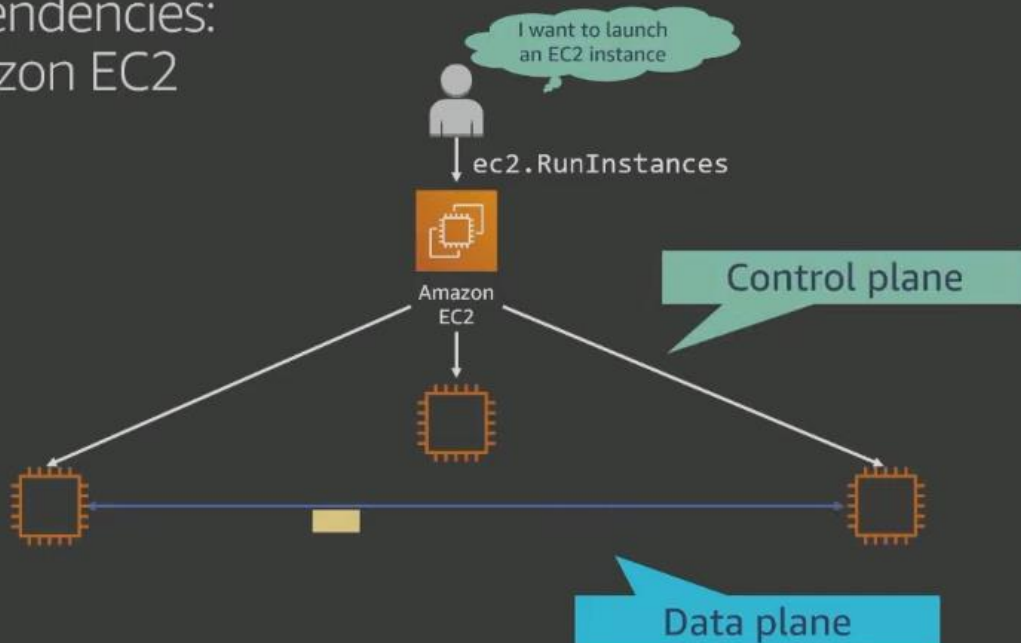
$$p * q < q$$



Not math:

Dependencies reduce availability

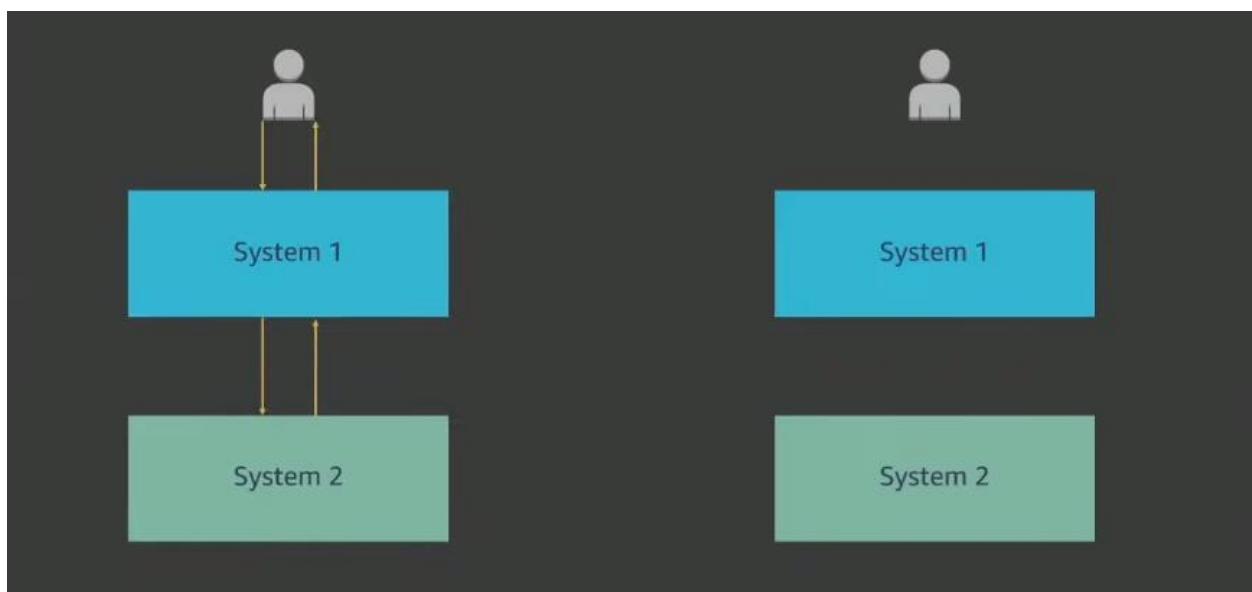
Dependencies: Amazon EC2



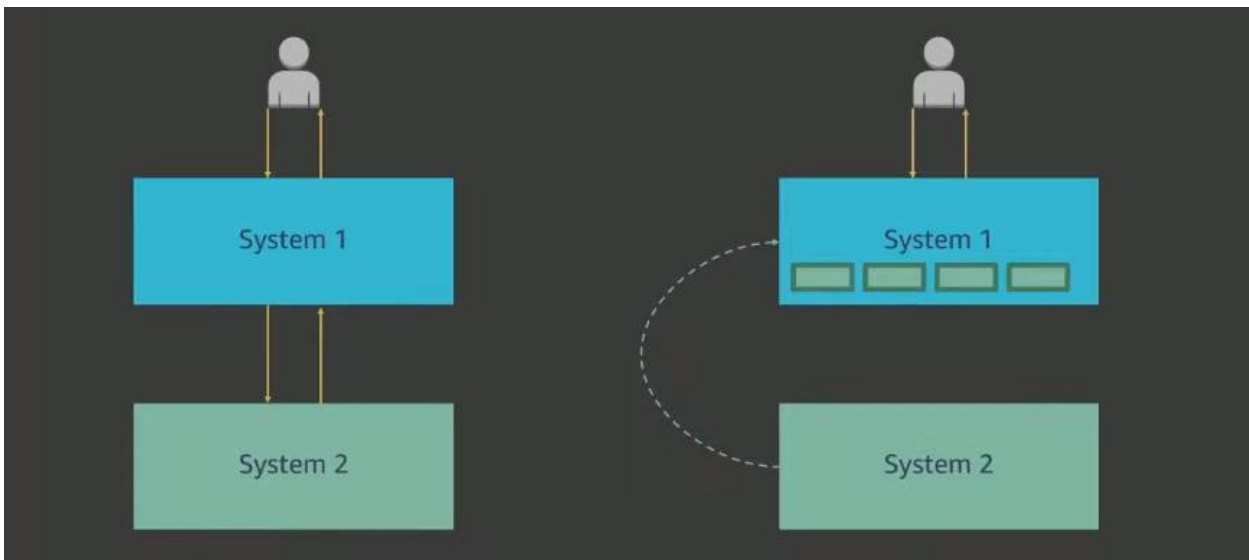
You might have 2 EC2 VMs running your apps with dependencies among them. We need to set up your Control plane for all these to happen for your Data plane to work.

Question 1: Which of these systems is more important to the success of your customers?

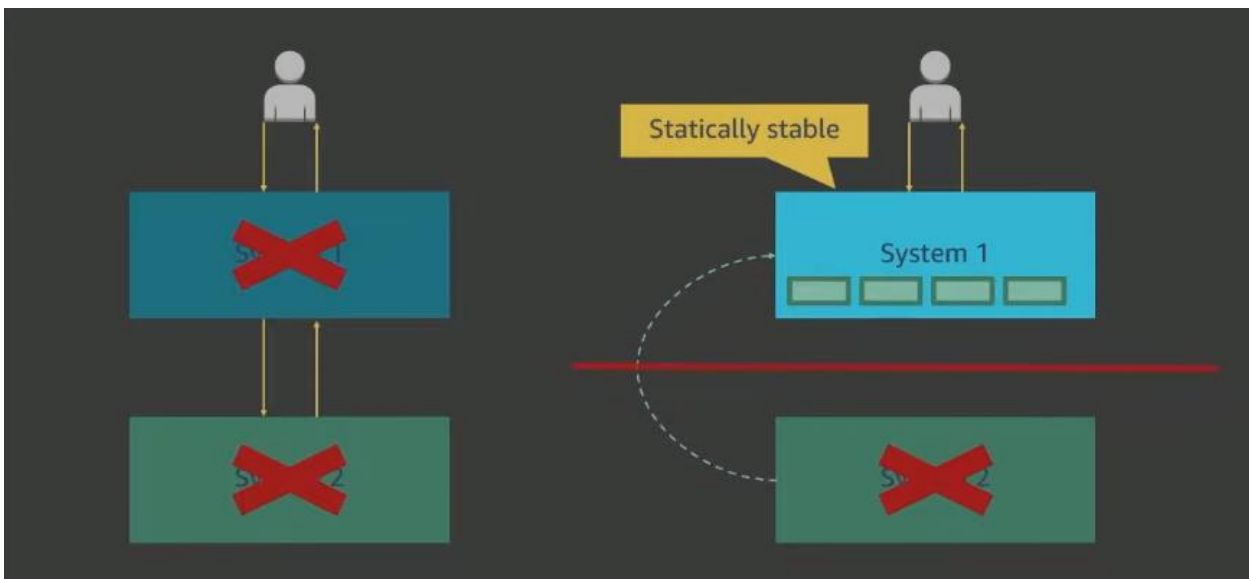
Question 2: Which of these systems has higher complexity?



There is a pattern in AWS for being able to use the data from another system but also being resilient in the case of data unavailability of system 2.

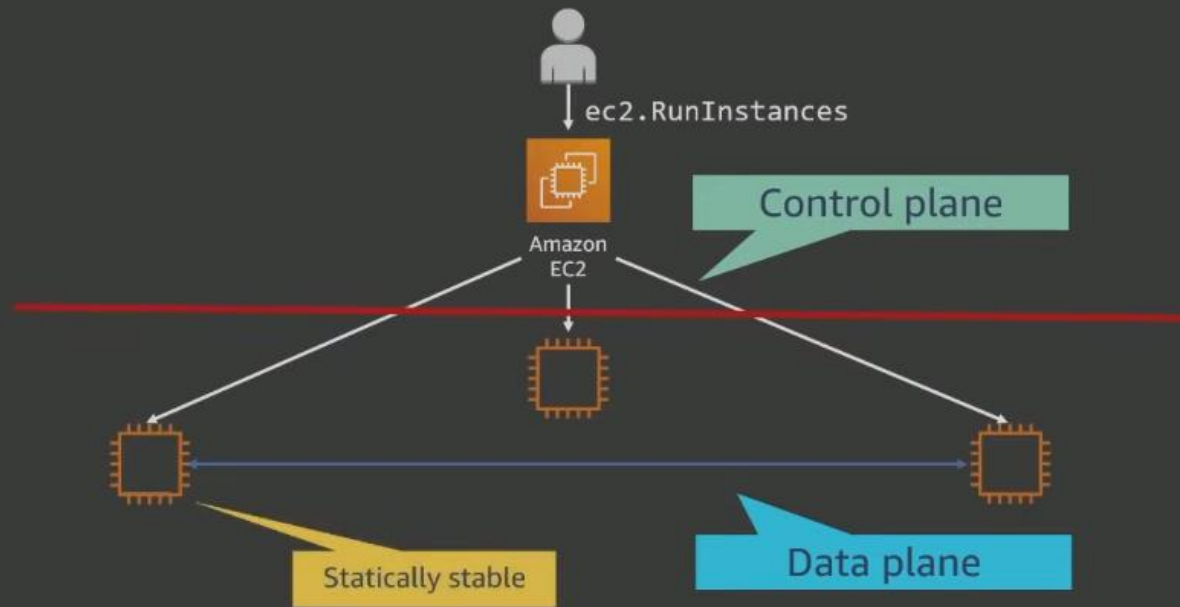


You can get the data into system 1 asynchronously so that it becomes available locally within system 1.



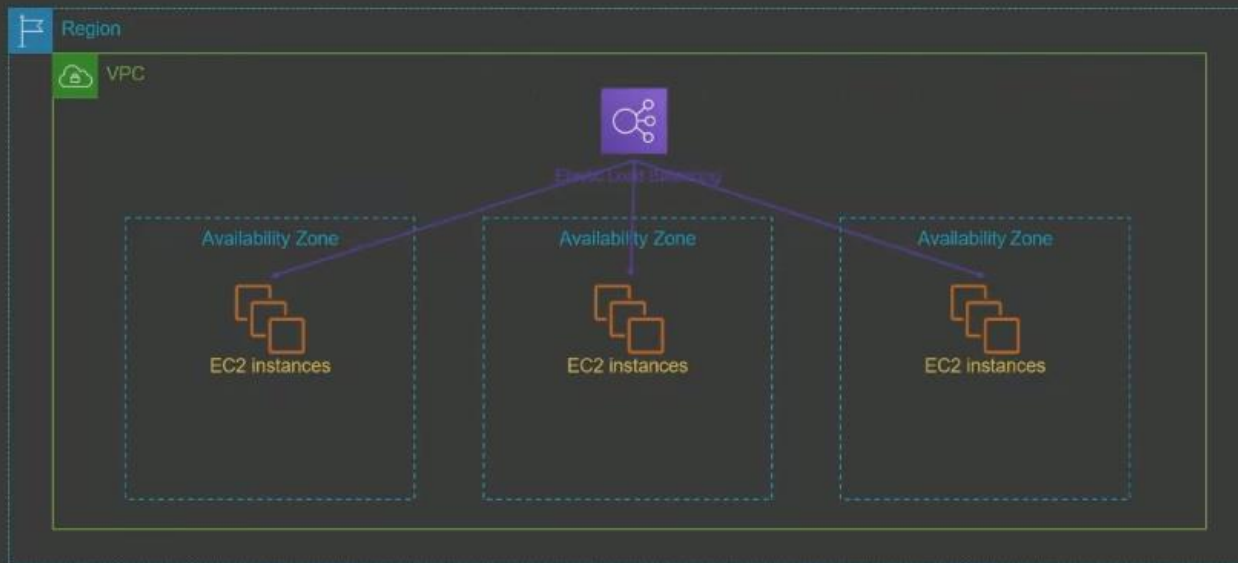
Now system 1 is still working albeit not updated

Static stability in practice: Amazon EC2

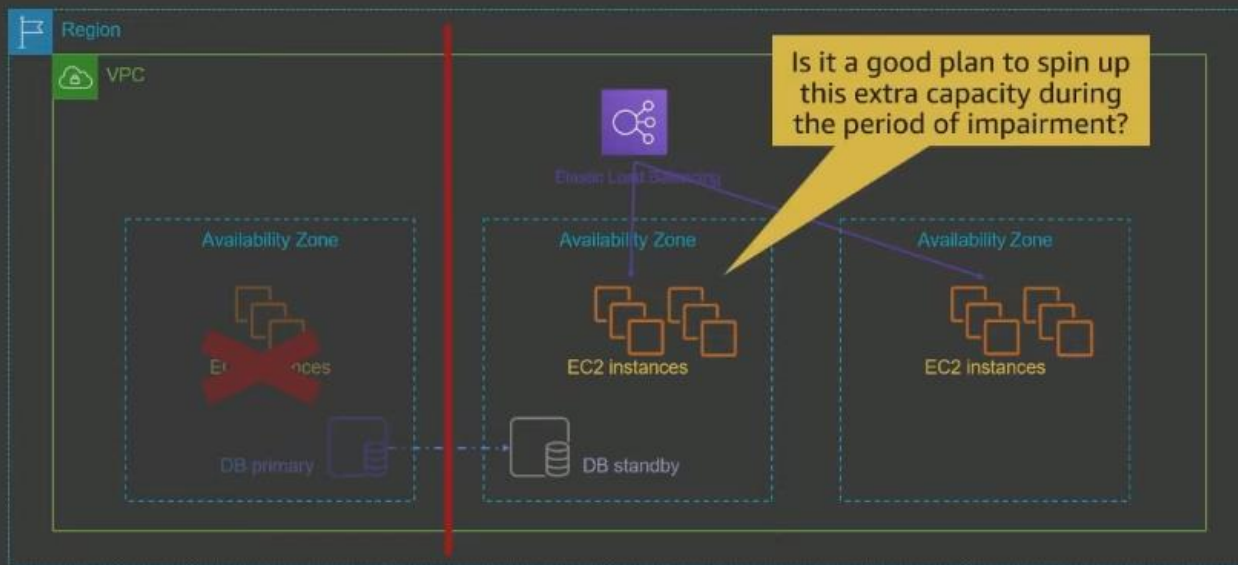


The Data plane is built to be statically stable, so that the Data plane doesn't even know about the existence of the Control plane.

Static stability in practice: Your high-availability workloads



Static stability in practice: Your high-availability workloads

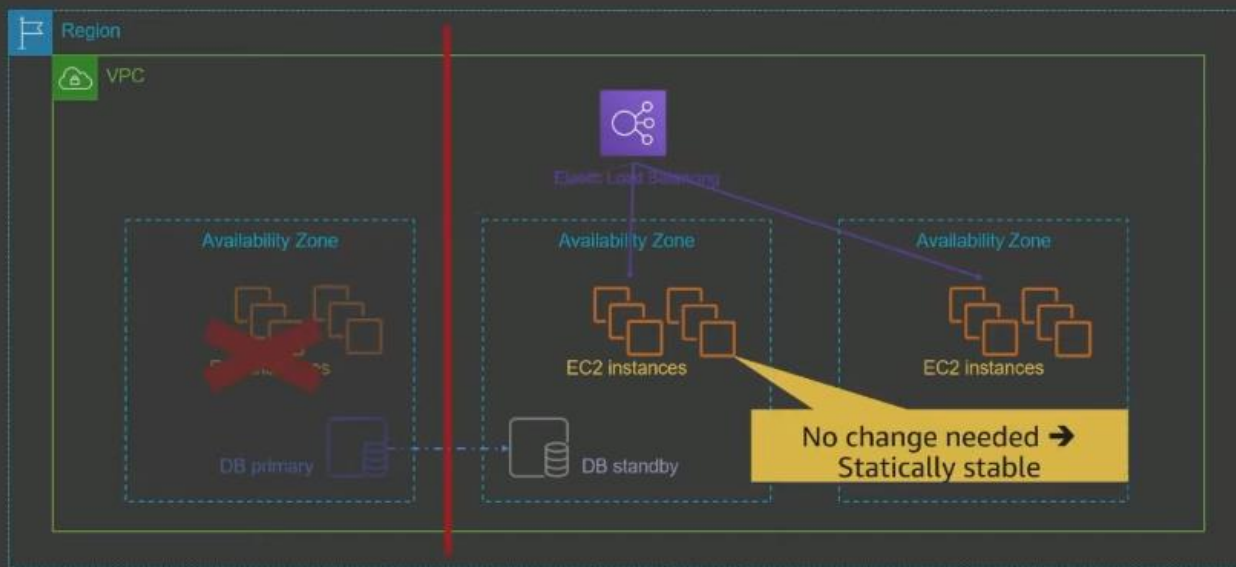


This is a rare event but possible.



We recommend that you not respond to an incident while it is raining but you should already be ready and planned for failure before it happens.

Static stability in practice: Your high-availability workloads



Static stability in practice: Your high-availability workloads



Static stability: Takeaways

- Consider the availability targets of your dependencies
- Use static stability techniques to avoid changing behavior during an impairment of your dependency
- Dependency \neq Destiny

Shuffle-sharding

This is a technique we use to isolate customer workloads from each other, this is one of the core competency of a data center management system that serves multi-tenants.

Private data center land



Shuffle sharding



Assume we have some web servers doing work

Shuffle sharding

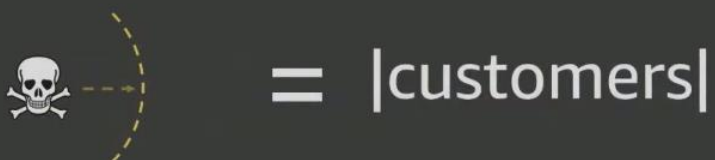


Assume there is a bad all within a server that shuts it down, the workload gets transferred to another server that gets shut down too

Shuffle sharding



Shuffle sharding



Shuffle sharding



You can do ordinary sharding by dividing into fleets and assigning customer workloads to each shards

Shuffle sharding



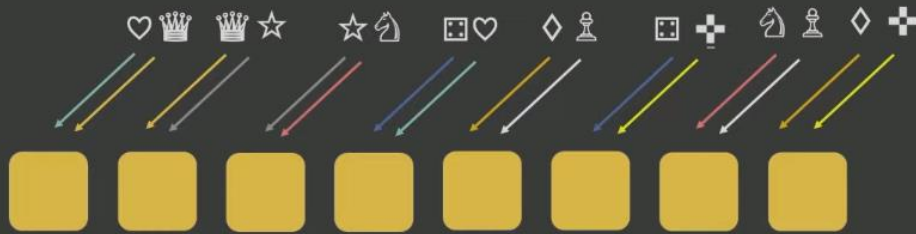
Now the other customers are isolated

Shuffle sharding

$$\text{skull icon} \rightarrow \text{bracket} = \frac{|\text{customers}|}{|\text{shards}|}$$

We can do more than this with the same number of resources using shuffle sharding, we use maths to assign customers smartly to a pair of server nodes

Shuffle sharding



Shuffle sharding



All you need is that the apps have some form of retry logic in them to allow the affected neighbors to get served elsewhere

Shuffle sharding

$$\text{skull icon} \approx \frac{\text{shardsize!}}{|\text{nodes}|!}$$

Shuffle sharding

Nodes = 8
Shard size = 2



Overlap	Probability
0	53.6%
1	42.8%
2	3.6%

Shuffle sharding



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

Shuffle sharding



Nodes = 8

Shard size = 2

Overlap	Probability
0	53.6%
1	42.8%
2	3.6%

Shuffle sharding



Nodes = 100

Shard size = 5

Overlap	Probability
0	77%
1	21%
2	1.8%
3	0.06%
4	0.0006%
5	0.0000013%

Shuffle sharding

Needs a client that retries or is fault-tolerant

Works for servers, queues, and other resources

Needs a routing mechanism: Per-customer DNS names, pre-resource DNS names, or a shuffle sharding-aware router

Retry, backoff, and jitter

Failures happen

Resilience through redundancy

Another challenge is how your clients in a distributed system achieve resilience and HA, by running multiple copies that enable retries when calls to other systems fail

Retries

Need availability?



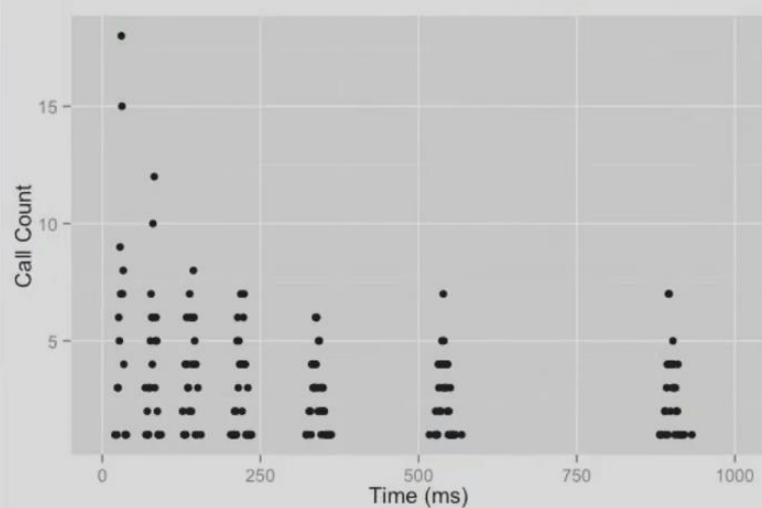
DoS your own service!

Retries when badly designed can lead to bad consequences

Thread#sleep()

We can defend against that by using exponential back-off retries in our clients or degrade to some stale data.

Distributed systems
have the
same problem



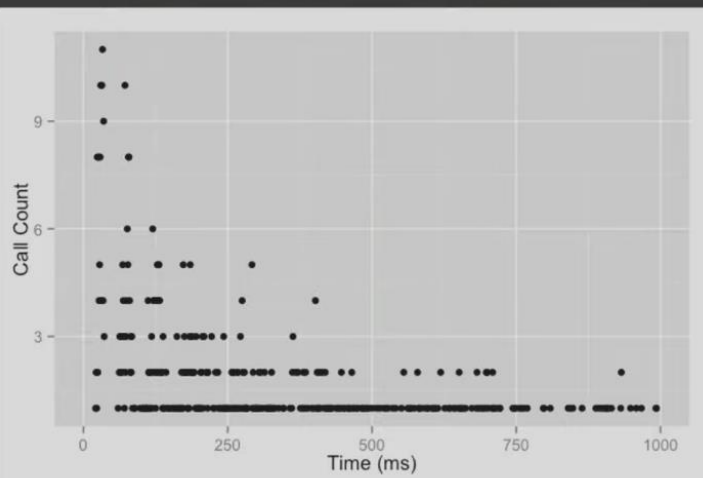
Humans tend to do things at specific times and this could still lead to DDOS or even cause less utilization of resources during down times

Jitter



Dados 4 a 20 caras.jpg: es:Usuario:Sabbur*derivative work: SharkD Talk CC BY-SA 3.0

Jitter allows us to introduce randomness into our systems



Doing things at random times is introducing jitter, while not requiring coordination between our systems

```
randint(0, base * 2 ** attempt)
```

This formula can be used to add jitter into our systems

How we do this at AWS

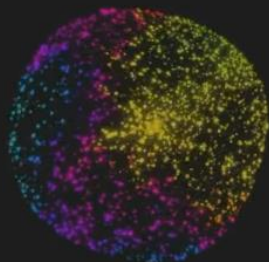
Always jitter when using backoff

Always jitter periodic work (timers, cron, etc.)

Consider adding jitter to all work

Meta lesson: Learn from each other

How does Amazon...



build resiliency into distributed systems?

approach DevOps?

engineer at scale?

The Amazon Builders' Library

How Amazon builds and operates software



Architecture, software delivery, and operations

By Amazon's senior technical executives and engineers

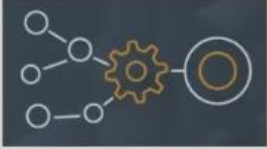
Real-world practices with detailed explanations

Content available for free on the website

The Amazon Builders' Library

ARCHITECTURE

LEVEL 300




Leader election in distributed systems

Author: Marc Brooker

Improving efficiency, reducing coordination, and simplifying architectures by using leader election.

SOFTWARE DELIVERY AND OPERATIONS

LEVEL 300




Going faster with continuous delivery

Author: Mark Mansour

Automating the software testing and deployment process for speed and reliability

SOFTWARE DELIVERY AND OPERATIONS

LEVEL 400




Implementing health checks

Author: David Yanacek

Automatically detecting and mitigating server failures without unintended consequences from fleet-wide false positives.

ARCHITECTURE

LEVEL 400



Workload isolation using shuffle-sharding

Author: Colm MacCarthaigh

Shuffle Sharding is one of our core techniques for drastically limiting the scope of impact of operational issues

Many more at aws.amazon.com/builders-library

Learn DevOps with AWS Training and Certification

Resources created by the experts at AWS to propel your organization and career forward



Take free digital training to learn best practices for developing, deploying, and maintaining applications



Classroom offerings, like DevOps Engineering on AWS, feature AWS expert instructors and hands-on activities



Validate expertise with the **AWS Certified DevOps Engineer - Professional** or **AWS Certified Developer - Associate** exams

Visit aws.amazon.com/training/path-developing/

© 2019, Amazon Web Services, Inc. or its affiliates. All rights reserved.

aws training and certification

Thank you!

Andrew Certain

@tacertain

Becky Weiss

becky@amazon.com

Colm MacCárthaigh

@colmmacc

David Yanacek

@dyanacek

Marc Brooker

@MarcJBrooker

aws
re:Invent

© 2019, Amazon Web Services, Inc. or its affiliates. All rights reserved.

