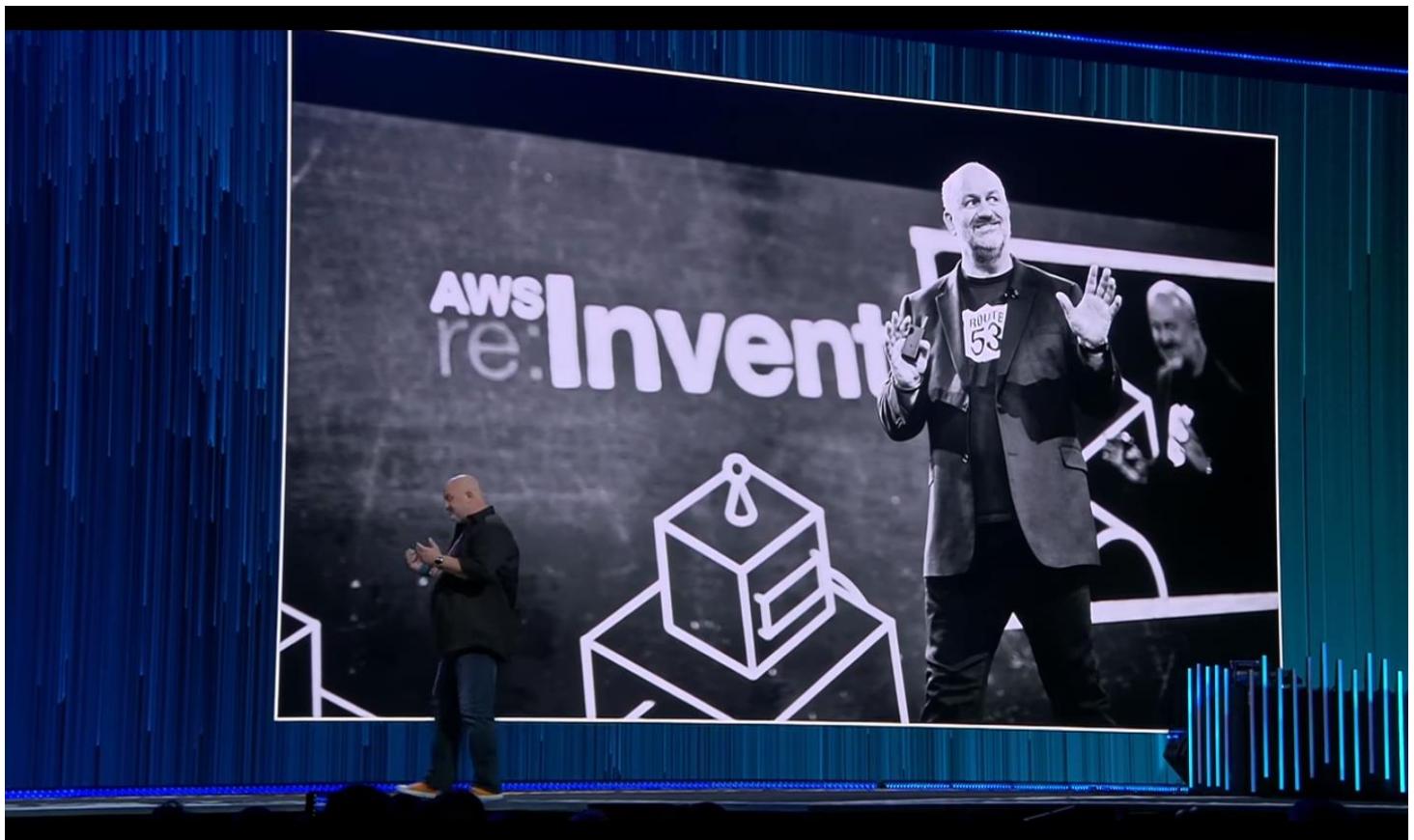
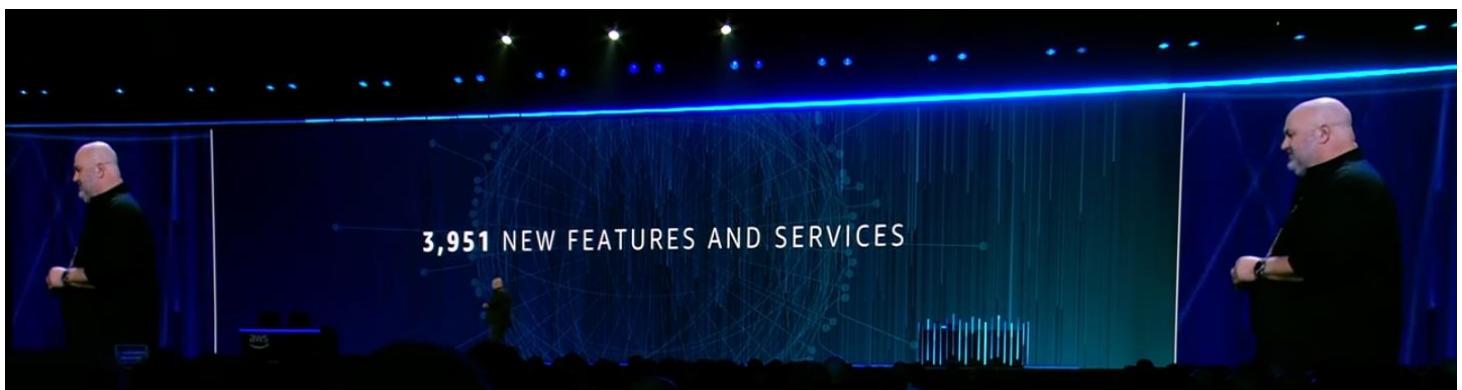
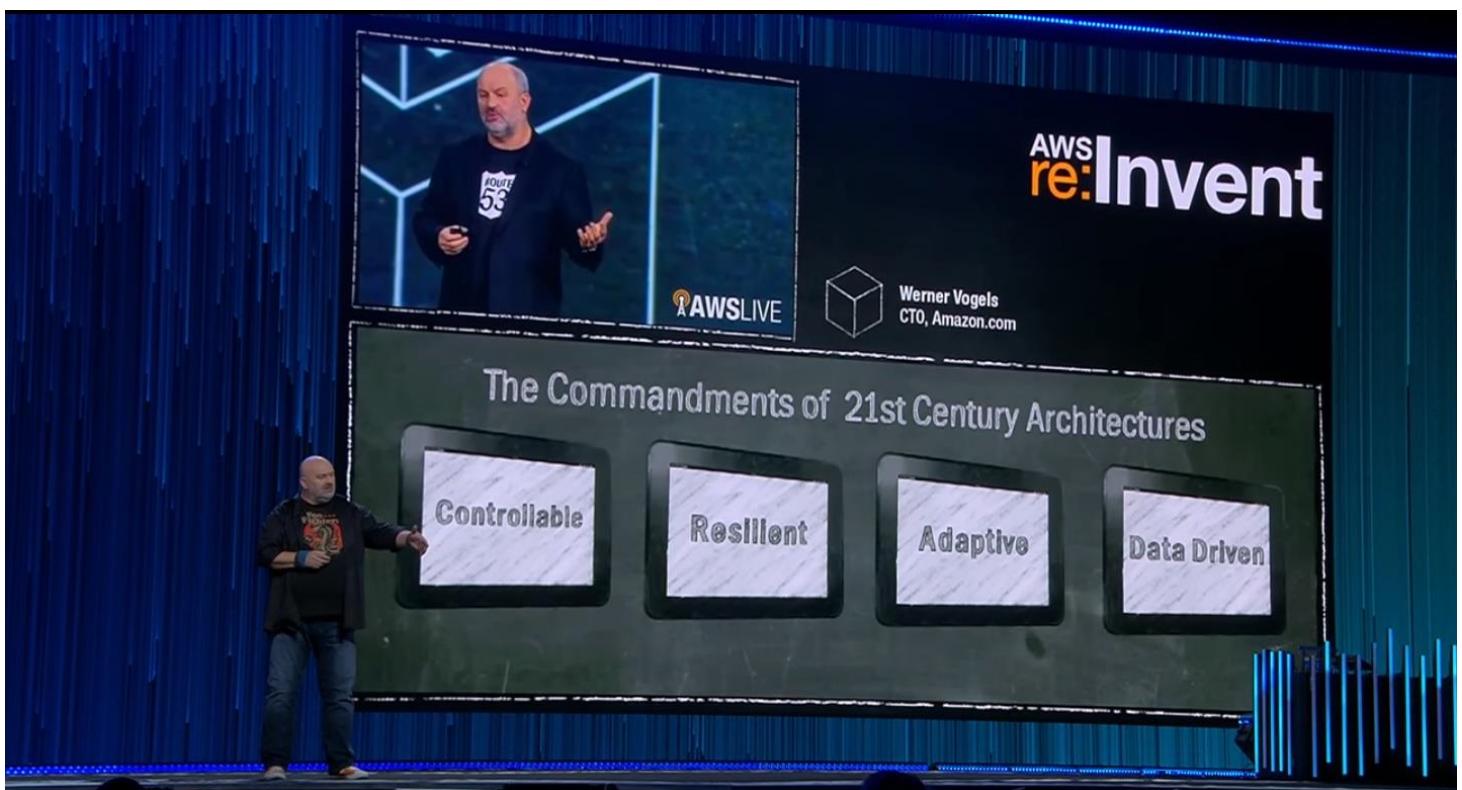




DR. WERNER VOGELS
CTO, AMAZON.COM







21ST CENTURY ARCHITECTURES

re:Imagined

21ST CENTURY ARCHITECTURES

COLLABORATIVE

TECHNOLOGY
DRIVERS TODAY



Data



IoT

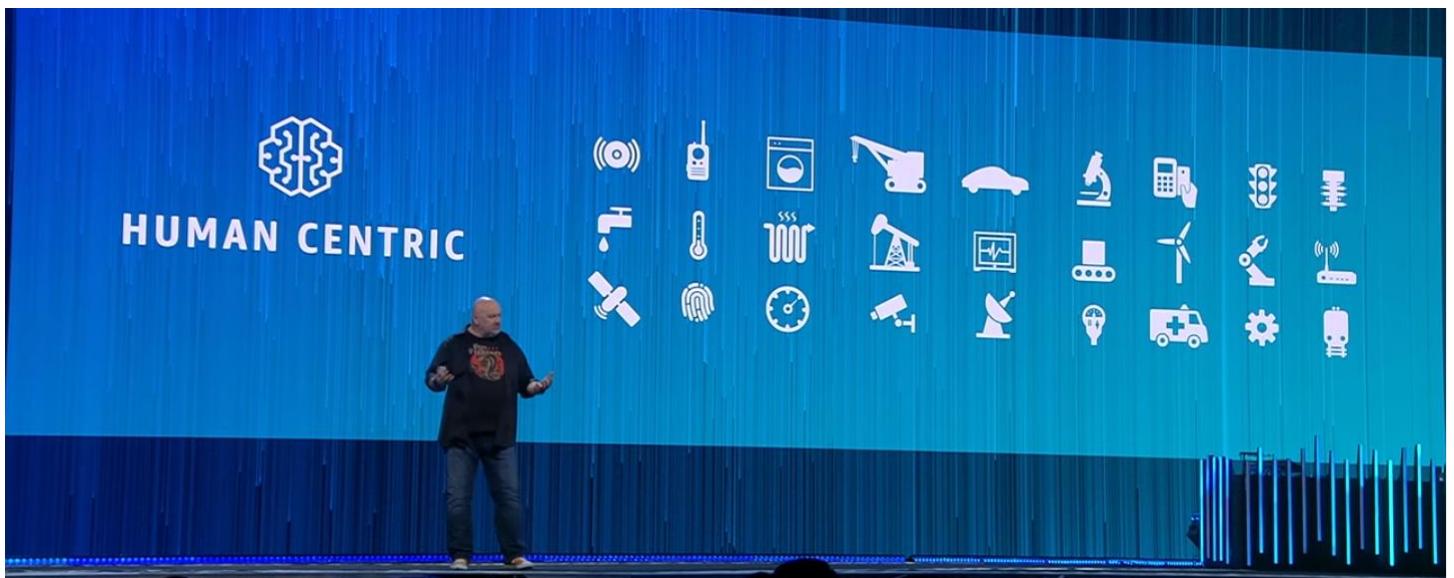
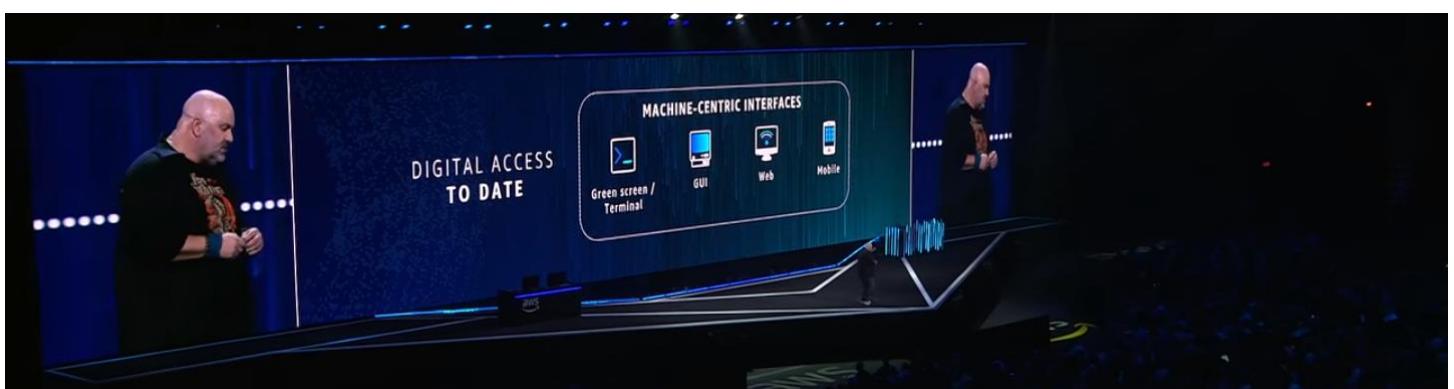


P3 instances



Deep Learning
Frameworks

We are now able to execute neural networks in real-time.





The first step to this is going to be via Voice, voice is the most natural interface for interacting in the world. We need to get this into our digital systems of the future.



The surgeons in the operating room should be able to use voice to interact with the machines around instead of having to take their hands off the patient. A sick father wants to talk to a device in natural language to get information about how to help his young child.



We needed to start building Voice Driven Systems.



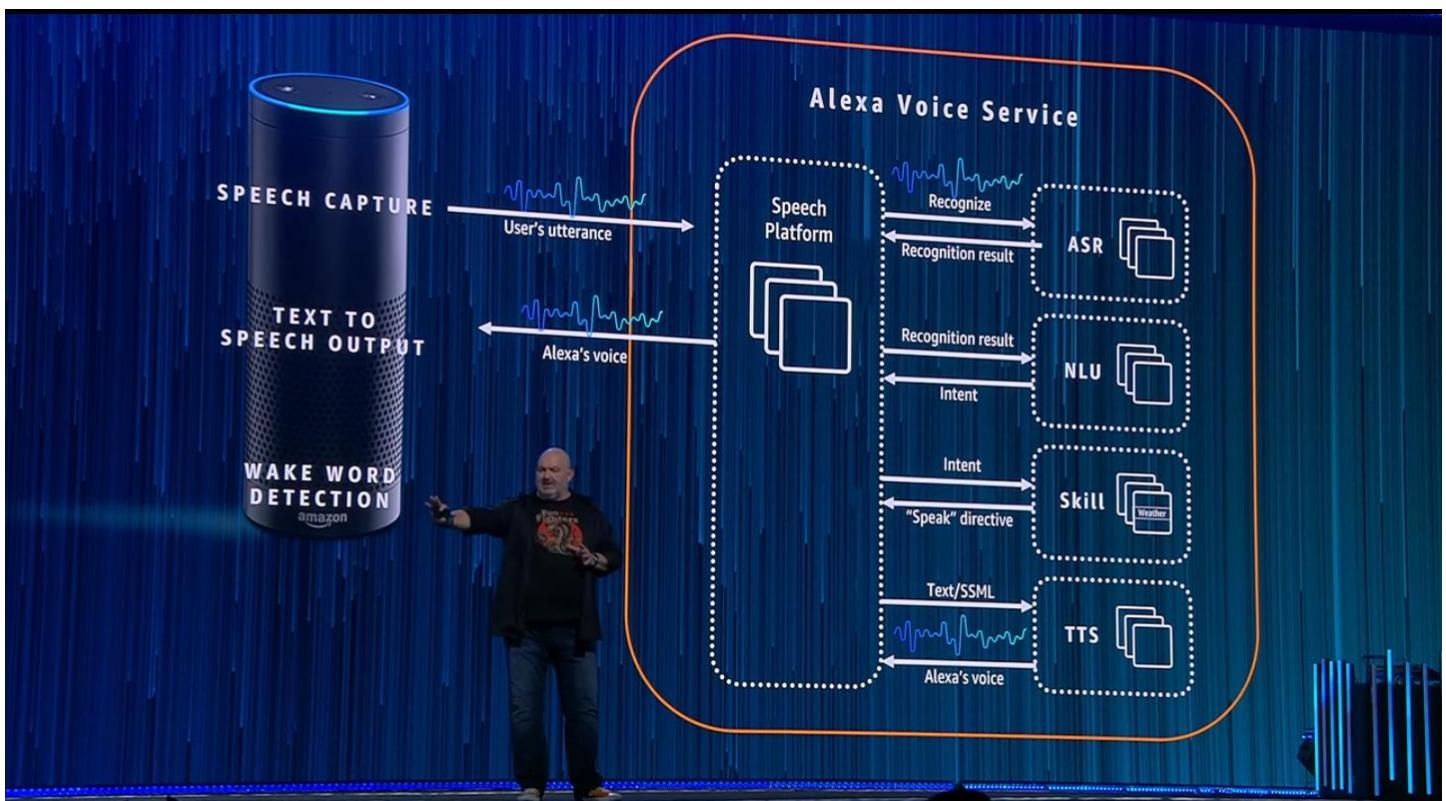
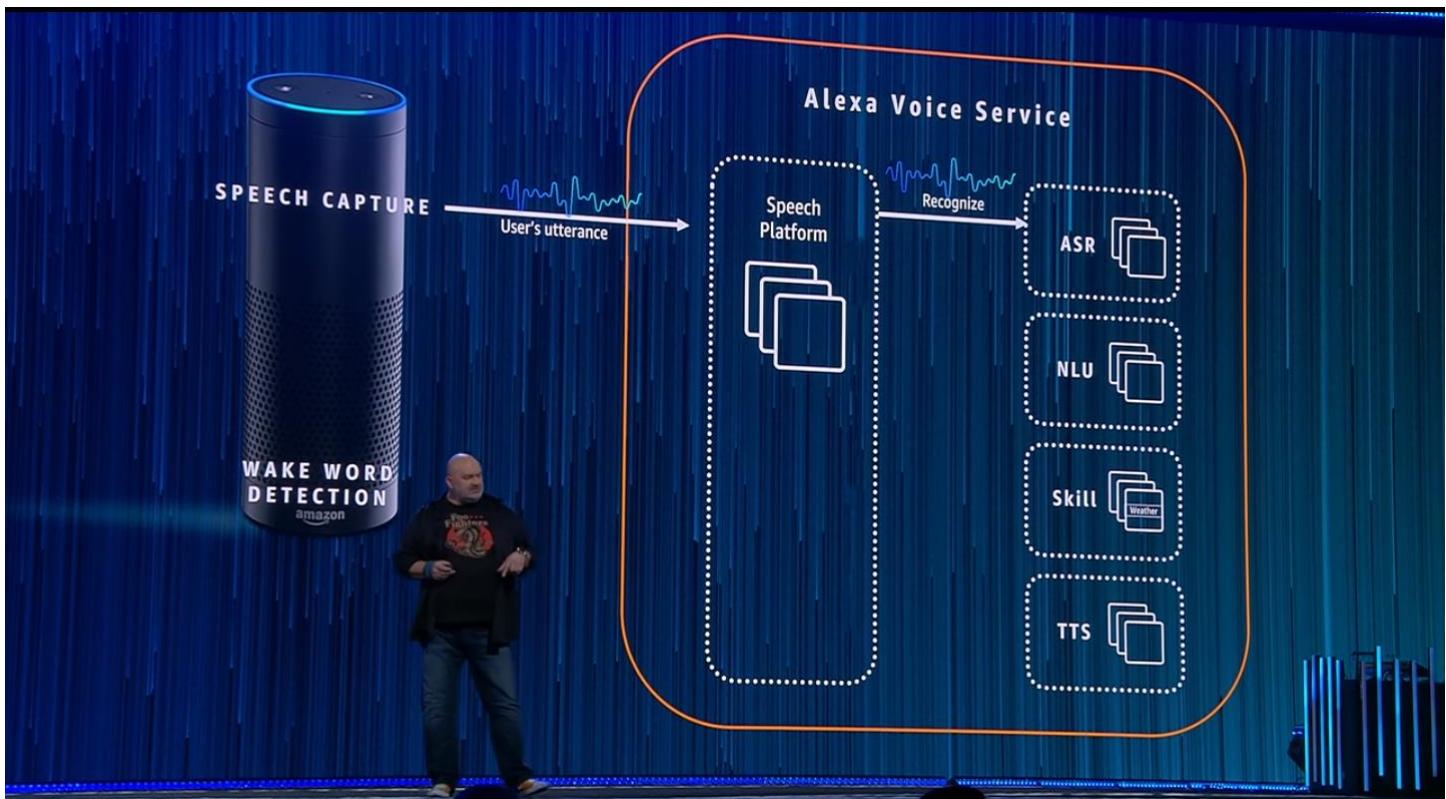
They put a voice interface to a digital system that a farmer can describe their plot of land and the machine learning system can query the digital system for information about what type and quantity of fertilizer to apply



All the smarts of these devices live in the cloud as the Alexa Voice Service and the Alexa Skills Kit.



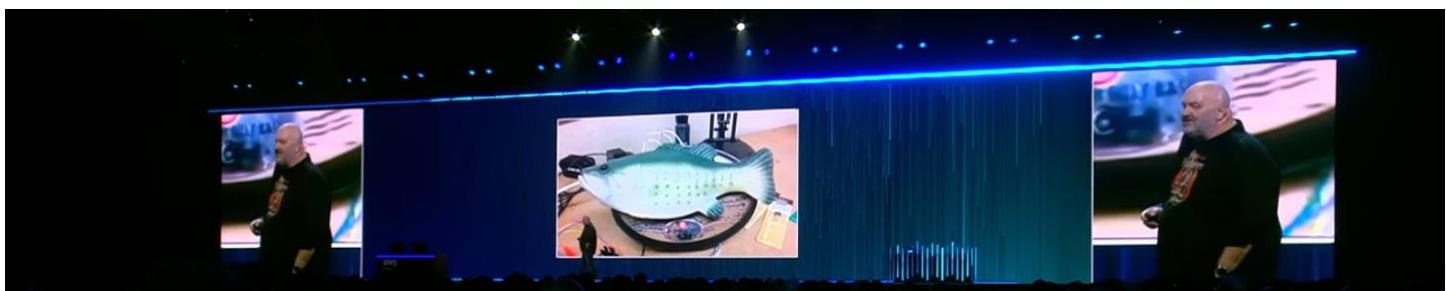
The wake word is 'Alexa', then the device starts recording and listening after that.

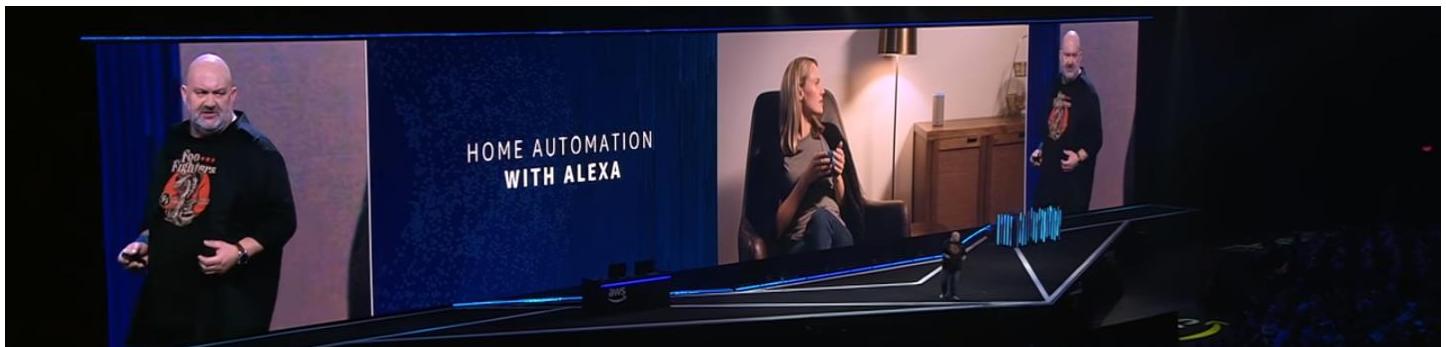


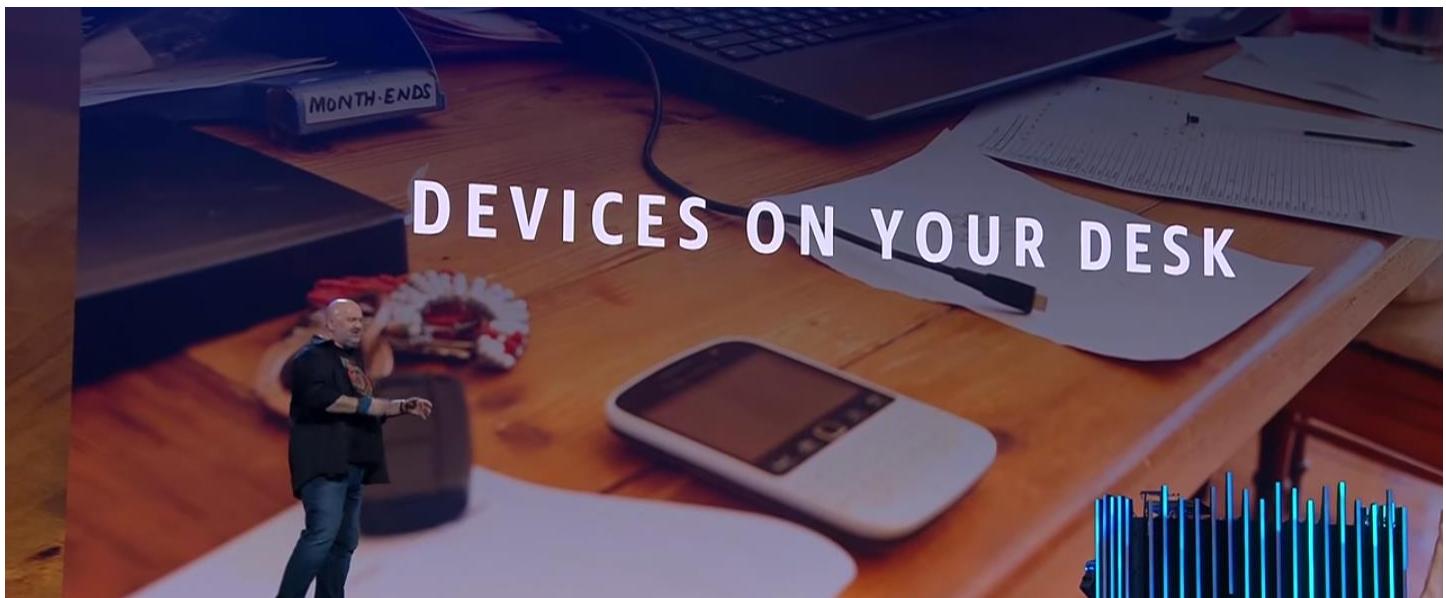
TTS creates the audio file that gets sent back to the device, which is then read out to the user. This is a pure cloud service. In AWS, Text-To-Speech (TTS) is the **Polly service**, the Natural Language Understanding (NLU) is the **Lex service**, you can build your own Alexa system if you want.



Anybody can integrate the Alexa Voice Service into their devices,







ALEXA HELPS YOU IN THE CONFERENCE ROOM

"Alexa, dim the lights."

"Alexa, lower the blinds."



ALEXA HELPS YOU IN THE CONFERENCE ROOM

"Alexa, ask Teem to
extend the meeting."

"Alexa, ask Teem to find an
open conference room."

teem

aws

ALEXA HELPS YOU AT YOUR DESK

"Alexa, join my meeting."

"Alexa, call Jeff."

aws



ALEXA HELPS YOU AT YOUR DESK

"Alexa, when is my next meeting?"

"Alexa, cancel my staff meeting."

"Alexa, schedule a meeting with Jeff."

aws







The **wynn hotel** is now putting an echo in each of their rooms so that guests can ask questions and get answers easily.



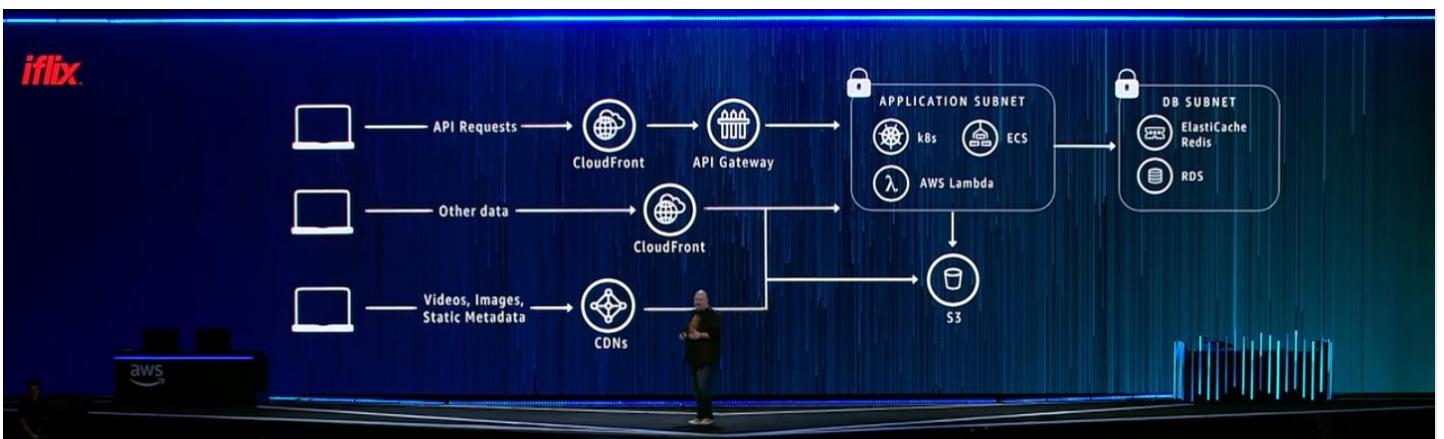
It will also mean that you are going to build your backend systems differently, the next generation of systems will be built using ***conversational interfaces*** because this is our systems will be winners in future.



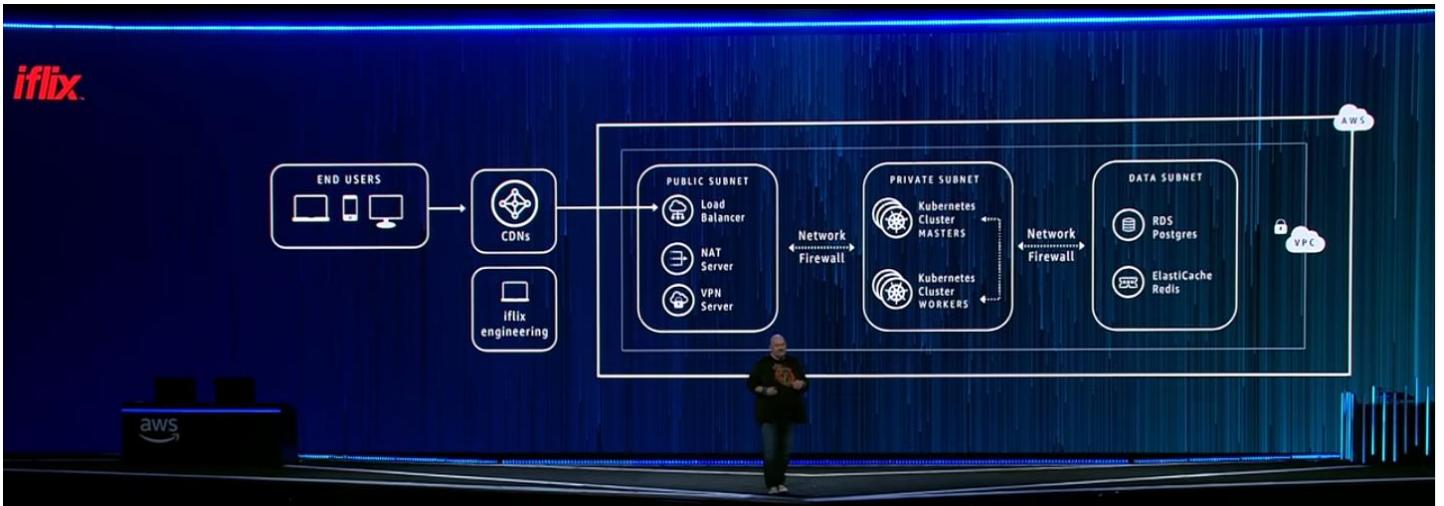
21st century architectures WILL have natural interfaces for voice, there are 3 different planes that systems developers need to be aware of



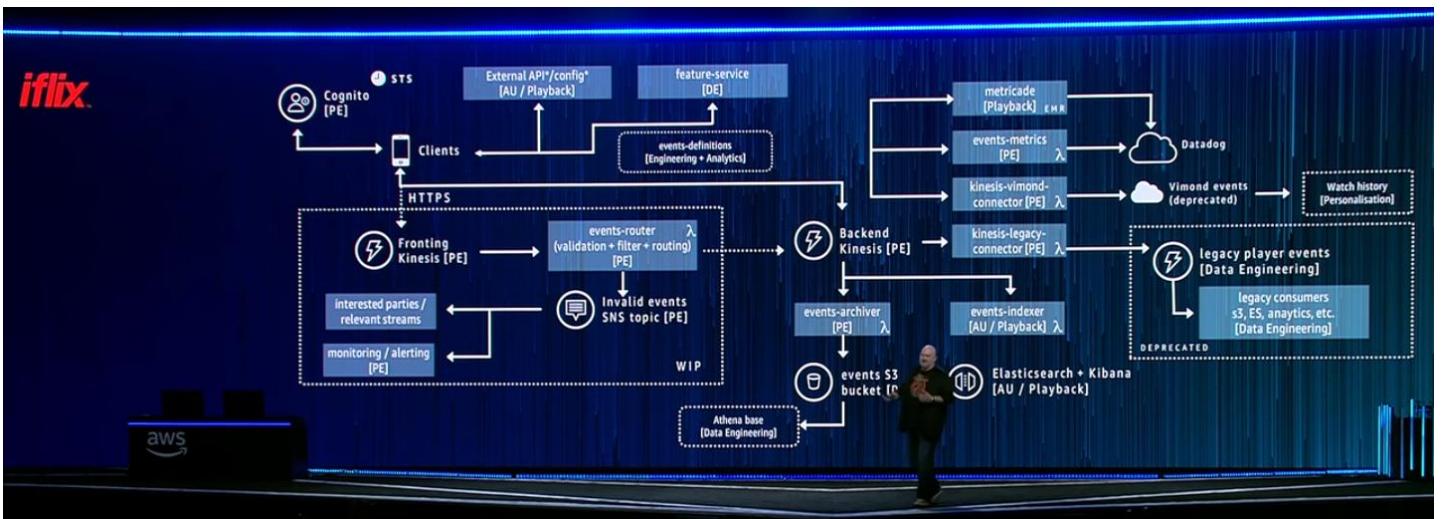
The Admin plane includes your deployment cycles, how you deploy, manage your containers and functions. The Control plane manages your deployed resources, the Data plane gives you access to your resources. There are different reliability and security requirements for these 3 planes.



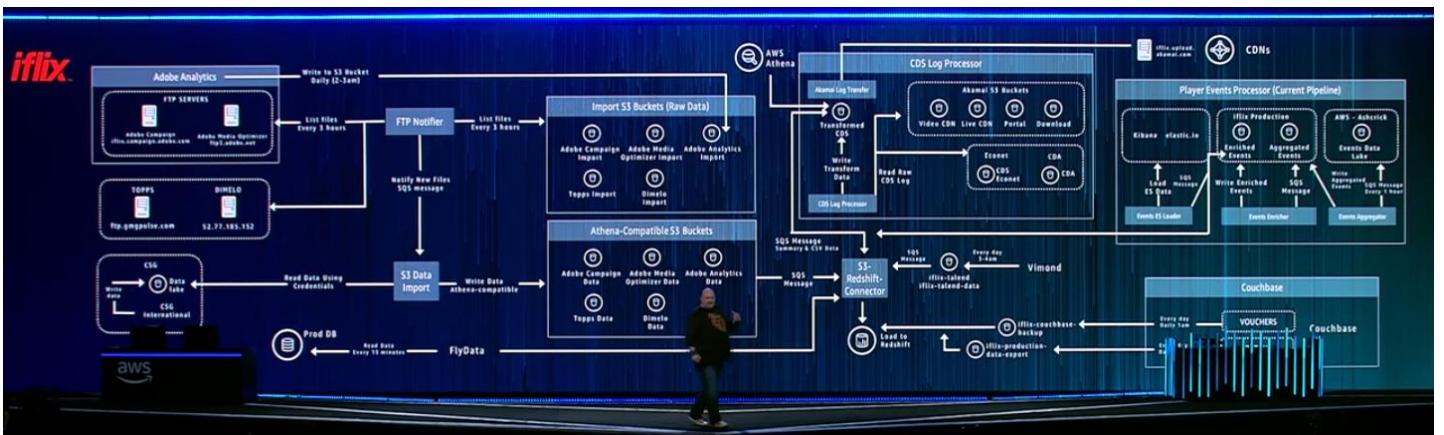
iflix is a video streaming company in EMEA, the typical technology architecture for a video streaming company is shown above but this is not all of it. Let's see what the rest looks like



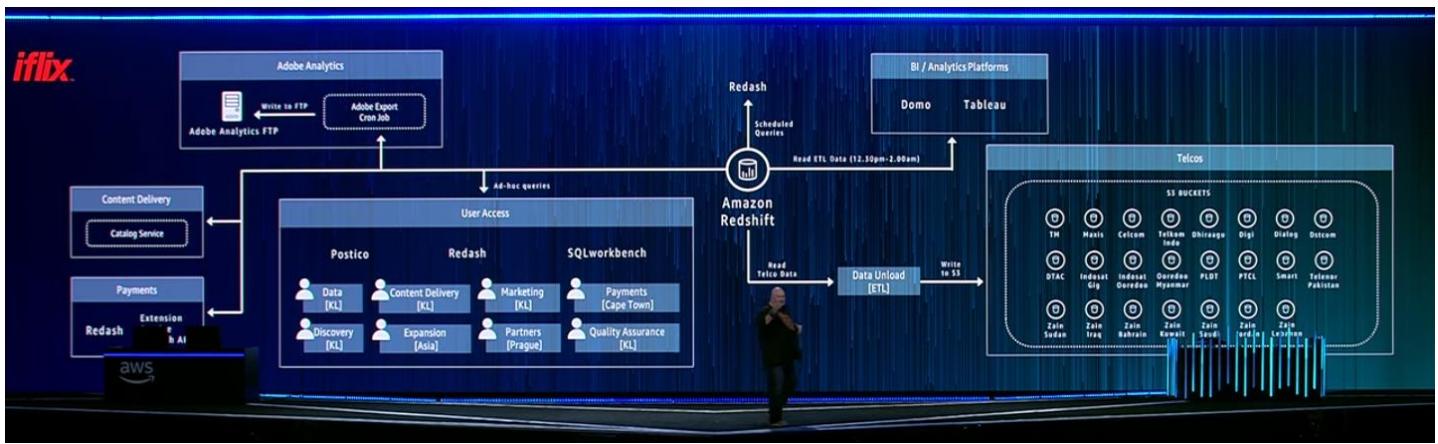
This is their **application layout** with subnets, VPC, k8s clusters and their associated databases running



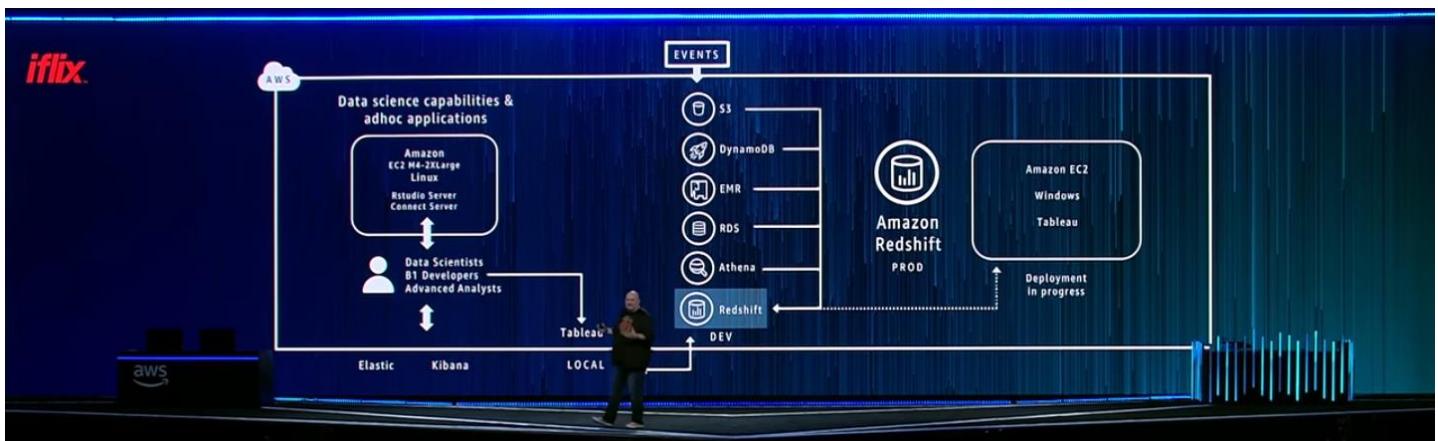
This is their **Eventing infrastructure/architecture** showing how events flow around their system



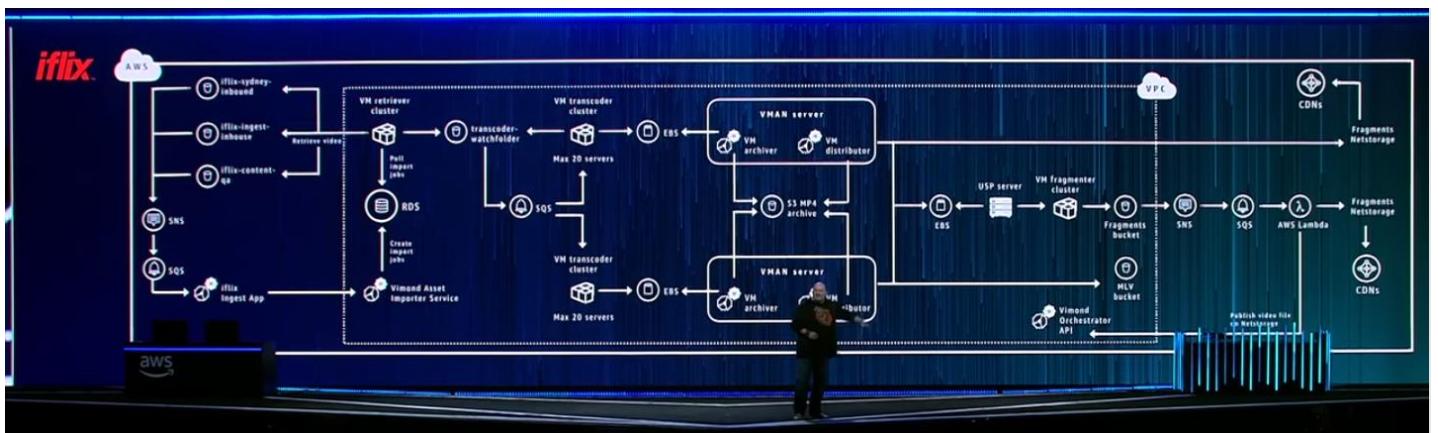
This shows where the **events** go



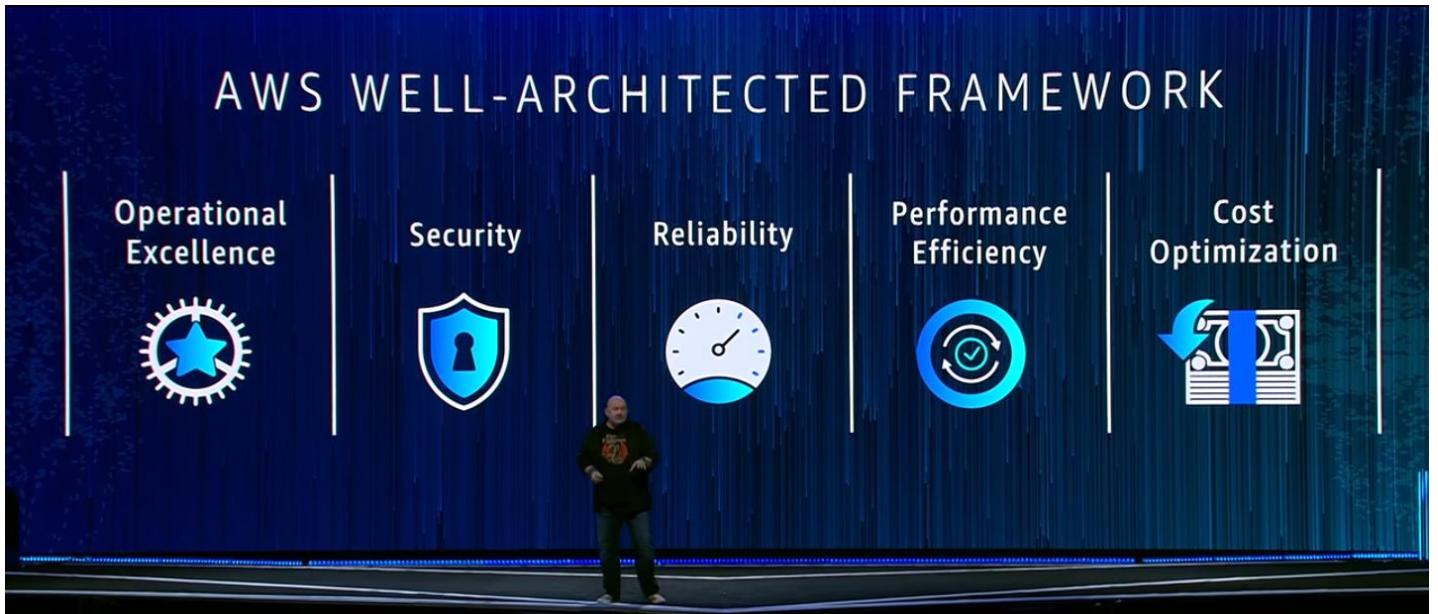
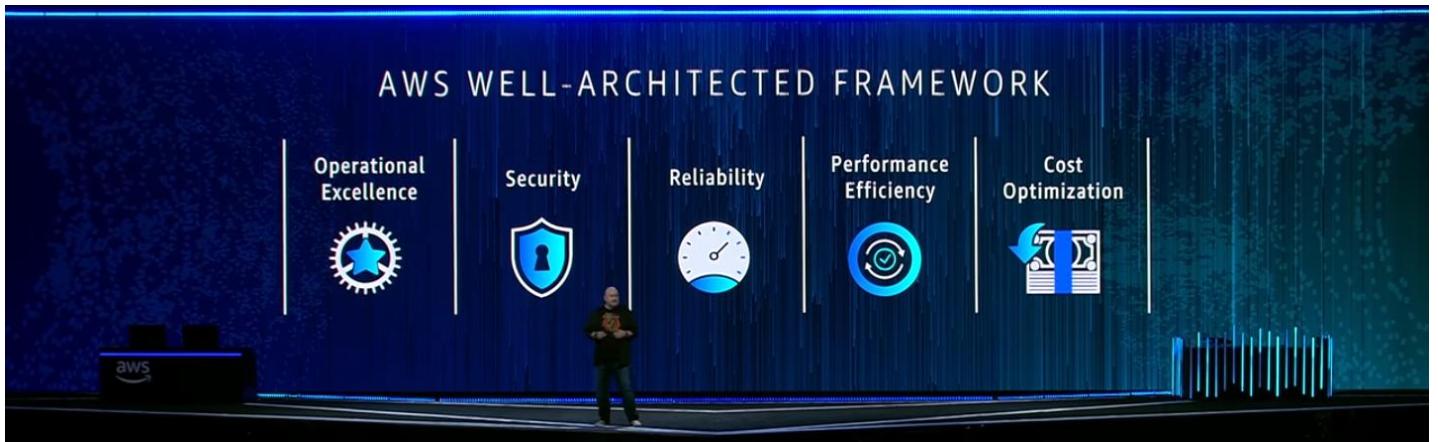
This is their Analytics infrastructure.



This is the environment where their ***data scientists can experiment*** with the data being collected to develop new products



This is their ***Video Data Ingestion environment*** about how to ingest video data in real time. AWS helps you to control complexity and keep things simple when building real architectures.



AWS Solution Architects have done extensive work in proposing 5 main pillars and give advice about to build very extensive architectures on top of AWS



Secure | https://aws.amazon.com/architecture/well-architected/

AWS Well-Architected

Learn, measure, and build using architectural best practices

AWS Architecture Center This Is My Architecture AWS Answers AWS Solutions Case Studies Cloud Security

AWS Well-Architected

The Well-Architected framework has been developed to help cloud architects build the most secure, high-performing, resilient, and efficient infrastructure possible for their applications. This framework provides a consistent approach for customers and partners to evaluate architectures, and provides guidance to help implement designs that will scale with your application needs over time.



Build and deploy faster
Stop guessing capacity needs, test systems at scale, and use automation to make experimentation easier by building cloud-native architectures.



Lower or mitigate risks
Understand where you have risks in your architecture, and address them before your applications are put into production.



Make informed decisions
Determine how architectural decisions and/or trade-offs might impact the performance and availability of your applications and business outcomes.



Learn AWS best practices
Access training and whitepapers that provide guidance based on what we have learned through reviewing thousands of customers' architectures on AWS.

Build using a structured approach

The AWS Well-Architected framework includes strategies to help you compare your workload against our best practices, and obtain guidance to produce stable and efficient systems so you can focus on functional requirements.

You can get started by downloading the AWS Well-Architected Framework whitepaper [PDF](#) | [Kindle](#)

Next, you can learn more about how to develop Well-Architected applications by downloading the whitepapers related to each of the five pillars of the Well-Architected framework.

<https://aws.amazon.com/architecture/well-architected/>



Make sure that you can evolve your systems over time and be able to be decomposed into smaller building blocks independently over time without having to change the overall system in one go.



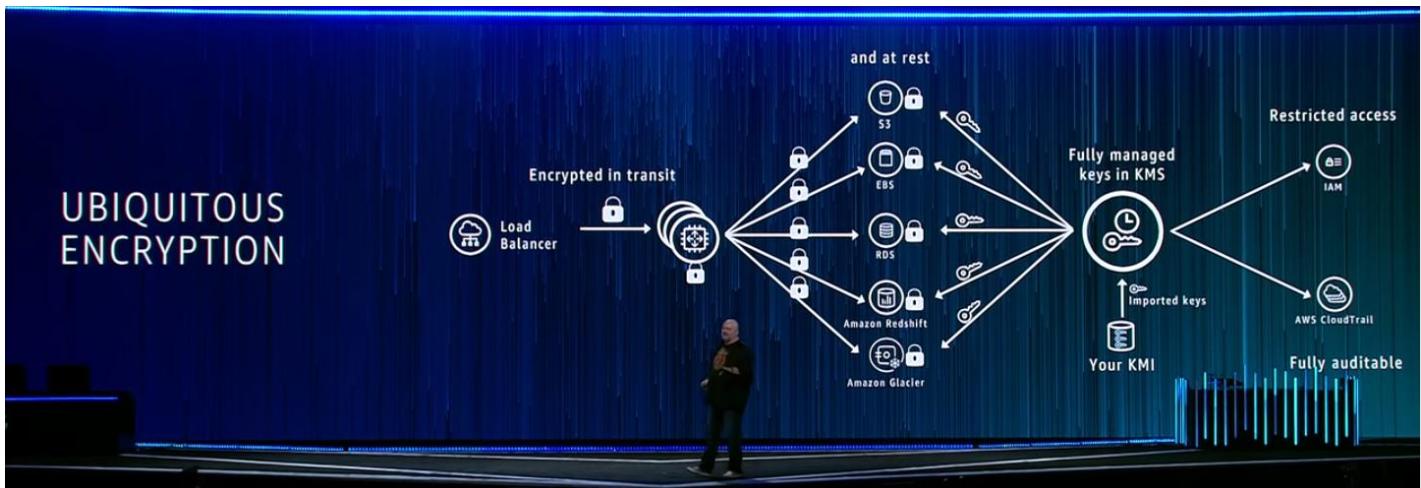
Protecting your customer and your business should be your number one priority, this comes before any feature development



Start taking away IAM roles until your engineers can no longer do their jobs, automate your best practices. You also need to have your playbooks and your runbooks ready every time



Encryption is the only way to know that you are the only one with access to your data. Signal-To-Noise toolkit can help you with this to encrypt data at rest and in transit

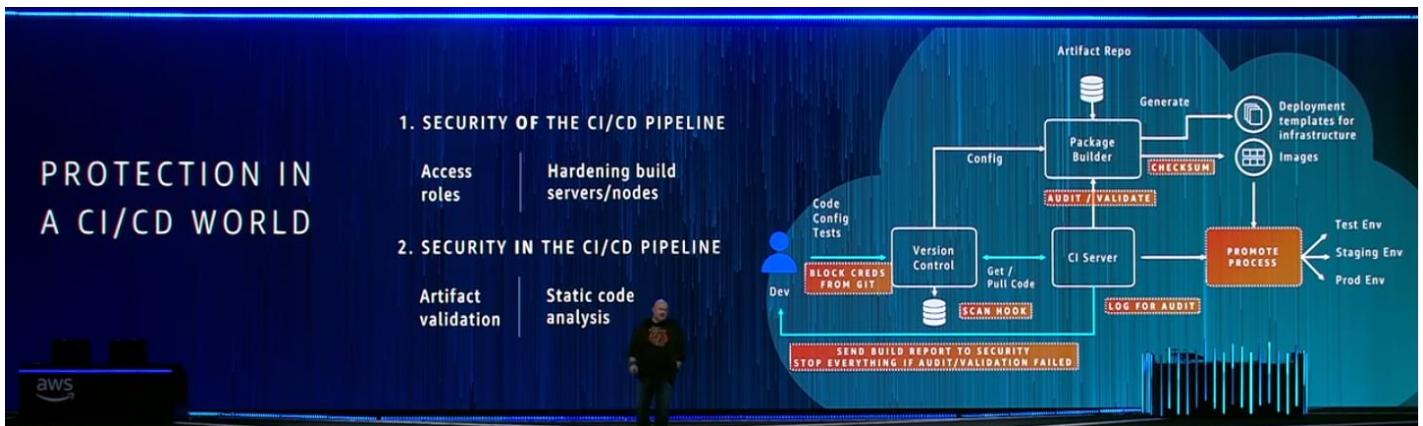
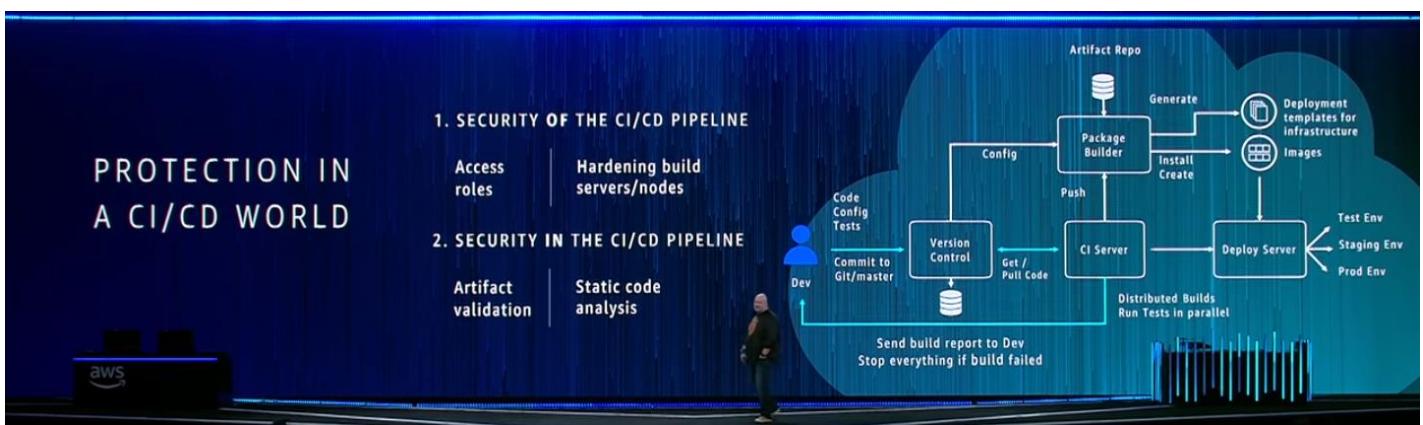
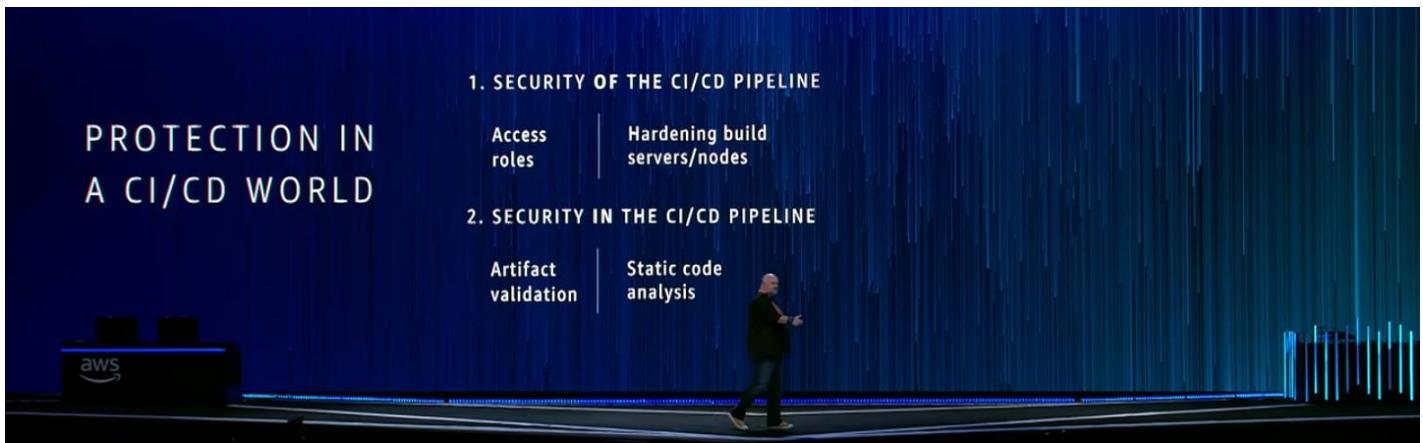


SECURITY IS EVERYONE'S JOB

Everyone should now be a security engineer, we all have to protect our customers

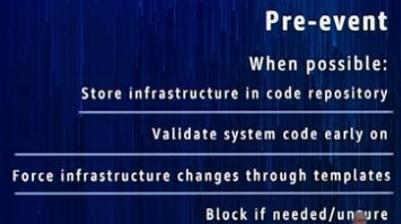


The best way to do this is via Automation



This is how your CI/CD pipeline should look like, each piece in the architecture should get security integrated into it to make sure that the software being developed is as secure as it can be.

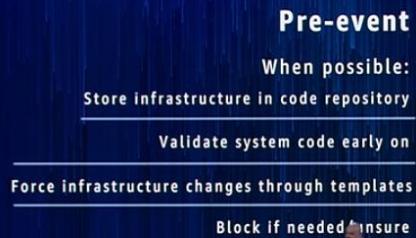
CONTROL AND VALIDATE



Post-event

- Always:
- Follow up on sensitive APIs
 - Use source of truth
 - Validate source
 - Decide on remediation

CONTROL AND VALIDATE



Triggers

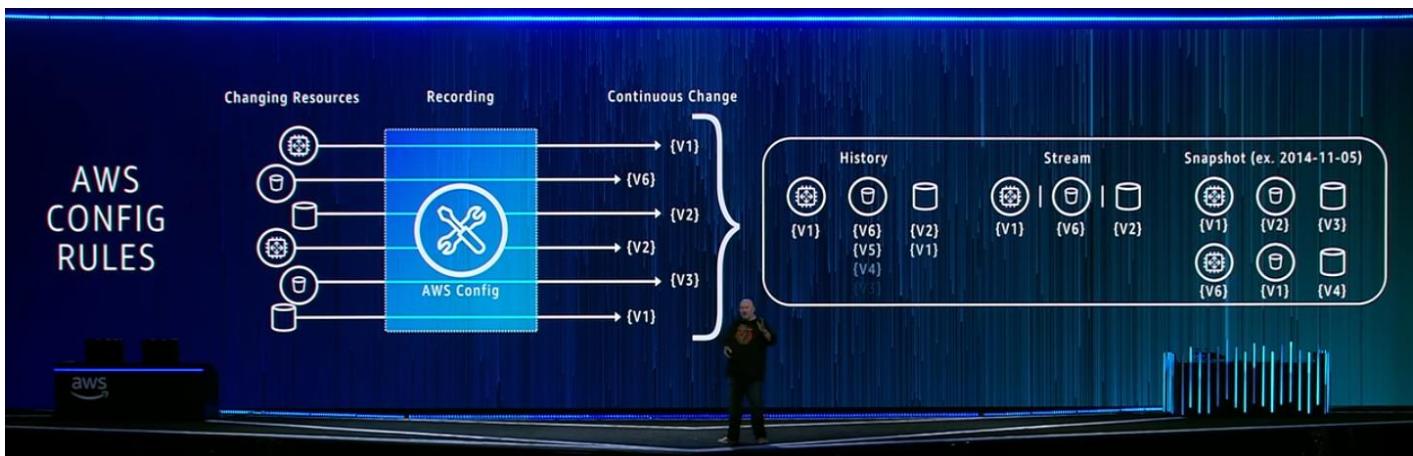
- Event based:
- Per change
 - Per day
 - Per framework

Any change to IAM needs to be manually inspected, every change should be audited

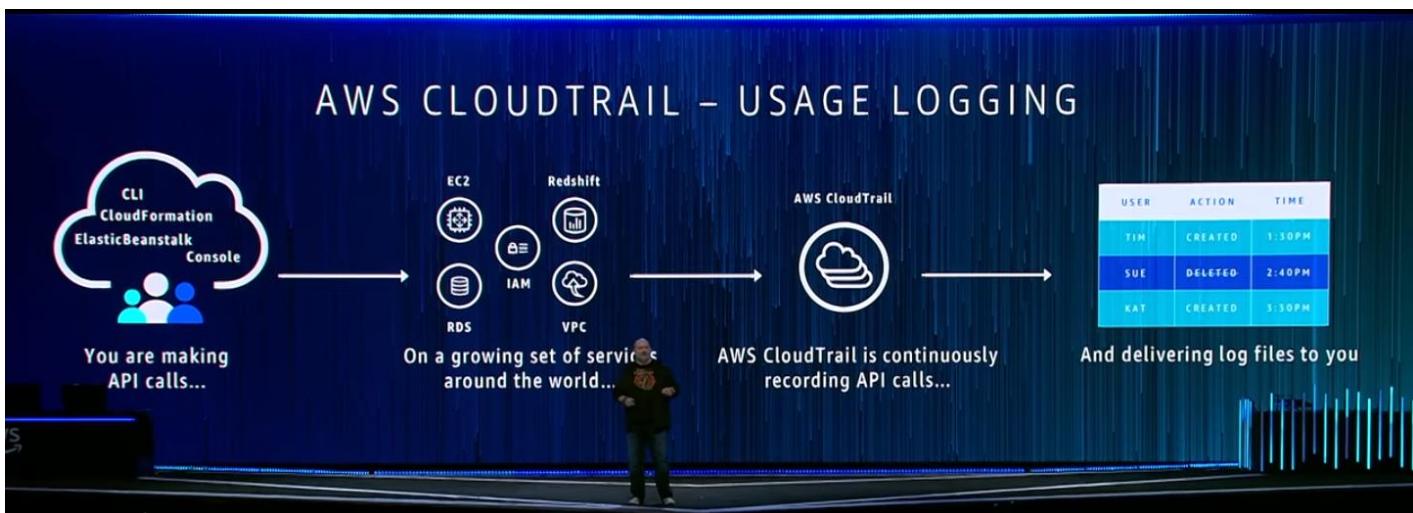
AUTOMATION HELP

- AWS Trusted Advisor
- AWS Config
- AWS Config Rules
- Amazon Inspector

- Amazon Cloudwatch
- AWS CloudTrail
- Amazon Macie
- Amazon GuardDuty



AWS Config Rules service can help you discover your whole environment, all the tools, services, and resources that you are using, and then it can start building an audit trail about all the changes that you are making to the resources that you are using. You can get streams as text, SNS notifications, create snapshots and compare them to how your system looked 2 weeks ago, etc.



AWS CloudTrail will log every API call to every one of your services, then put the logs into secure buckets in S3 for you to do analysis on who is accessing your systems and put alarms on them



EVERY GREAT PLATFORM HAS A GREAT IDE

aws



AWS Cloud9

A cloud IDE for writing, running,
and debugging code

Generally available today

aws



Zero to
deploy FAST



Deep
integration



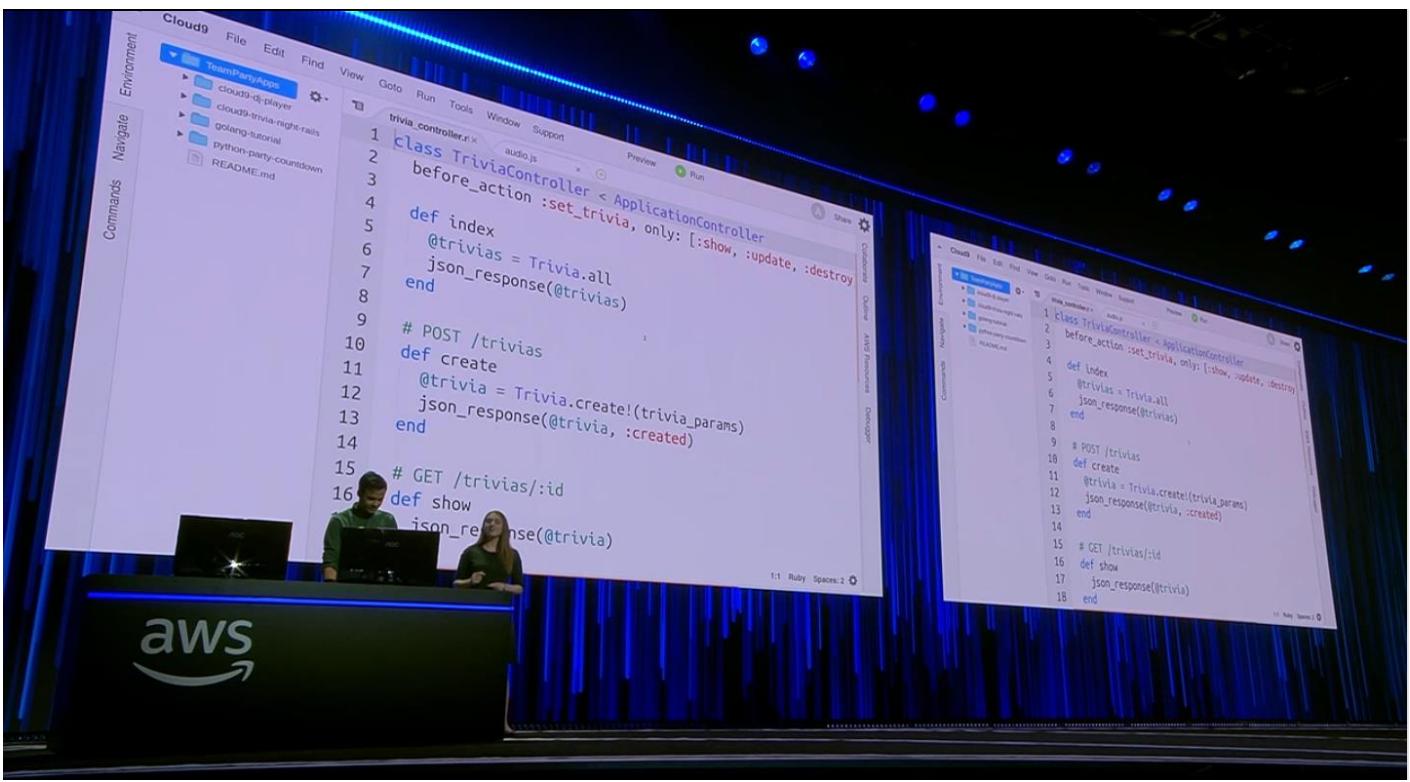
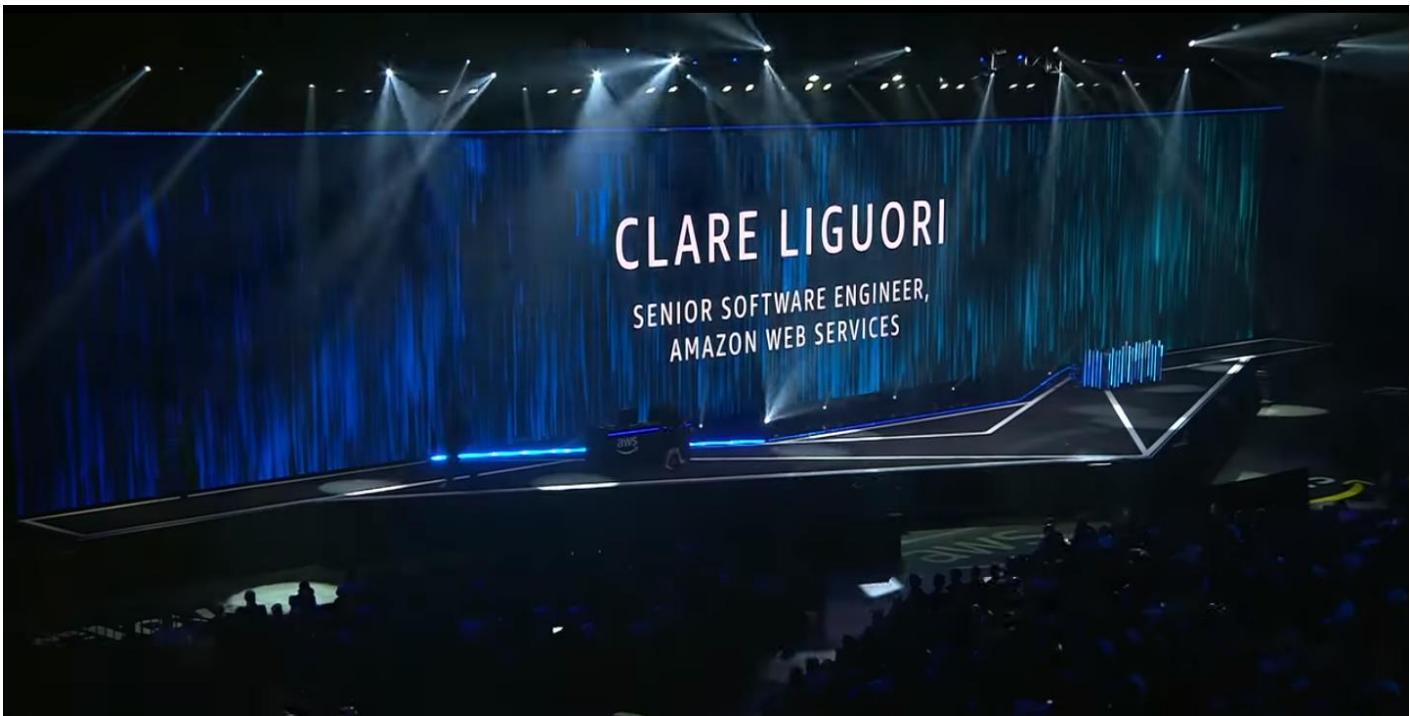
Real-time pair
programming



Broad debugging
support

aws

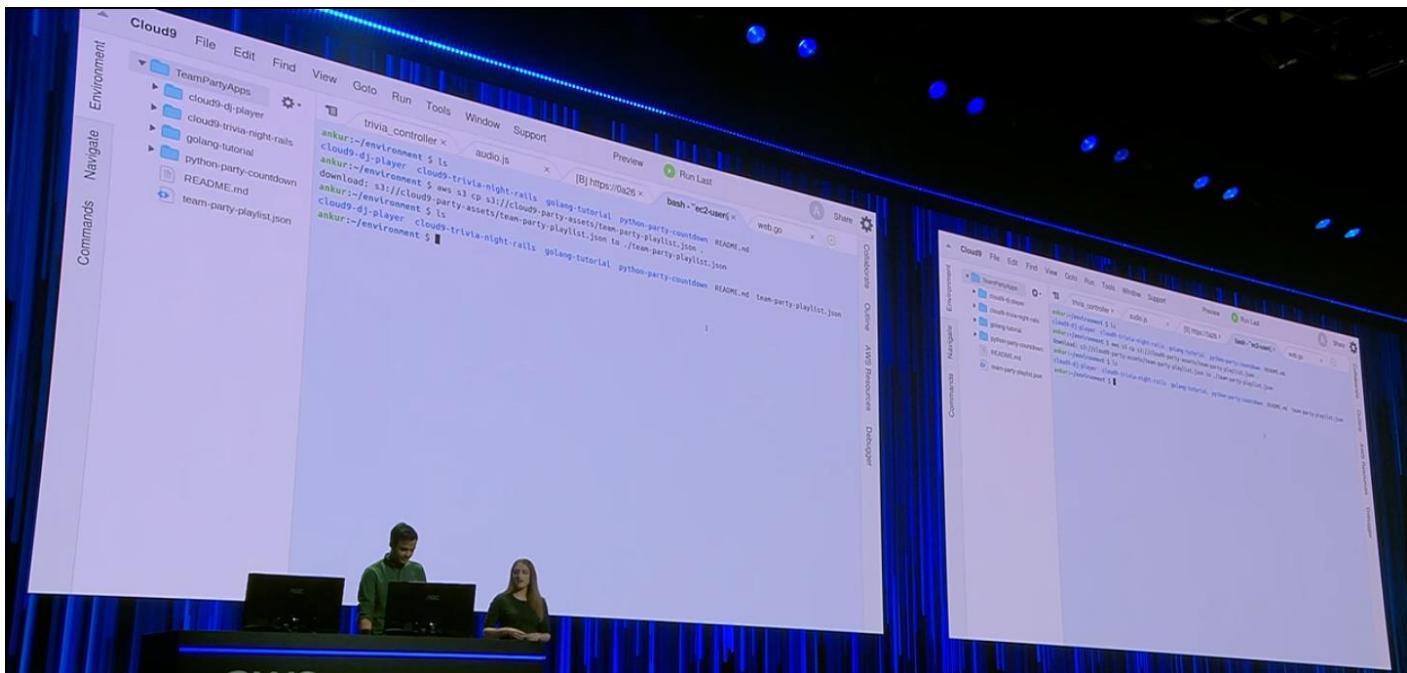
Includes the ability to debug your lambda functions in the cloud



This is Ankur's Cloud9 environment that we are able to pull up on any browser from anywhere,

A screenshot of the Cloud9 IDE interface. On the left, there's a sidebar with tabs for 'Commands', 'Navigate', and 'Environment'. The 'Environment' tab is selected, showing a file tree for 'TeamPartyApps' containing files like 'cloud9-dj-player', 'cloud9-trivia-night-rails', 'golang-tutorial', 'python-party-countdown', and 'README.md'. In the main workspace, there are two terminal tabs: 'trivia_controller.rl' and 'audio.js'. The 'audio.js' tab shows a command-line session on an EC2 user instance:

```
ankur:~/environment $ ls
cloud9-dj-player cloud9-trivia-night-rails golang-tutorial python-party-countdown README.md
ankur:~/environment $
```



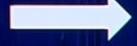
The AWS CLI comes pre-installed and can be used to access AWS resources using your automatically attached account credentials

A screenshot of the Cloud9 IDE interface. On the left, there's a sidebar with tabs for 'Commands', 'Navigate', and 'Environment'. The 'Environment' tab is selected, showing a file tree for 'TeamPartyApps' containing files like 'cloud9-dj-player', 'cloud9-trivia-night-rails', 'golang-tutorial', 'python-party-countdown', and 'README.md'. In the main workspace, there are two terminal tabs: 'trivia_controller' and 'aud'. The 'trivia_controller' tab shows a command-line session:

```
36
37 function orderDetails() {
38     var orderItemData = [
39         {
40             "itemName": "Large Pepperoni Pizza",
41             "quantity": 2
42         },
43         {
44             "itemName": "Beer",
45             "quantity": 2
46         }
47     ];
48
49     var customer = {
50         name: 'Cloud9 Team',
51         address: new pizzapi.Address('MGM Grand'),
52         phone: '206-555-5555',
53     }
```

Below the terminal, there's a video feed of two people, a code editor with 'index.js' open, and a group chat window where one person says 'Hey Ankur, I have one more thing to add.' and another responds 'Go for it!'. The right side of the screen shows a preview of a web application with a 'Run' button.

PUBLISH DIRECTLY INTO CODESTAR TOOLS



AVAILABILITY, RELIABILITY, AND RESILIENCE IN 21ST CENTURY ARCHITECTURES



AVAILABILITY,
RELIABILITY,
AND RESILIENCE
IN 21ST CENTURY ARCHITECTURES

Test recovery procedures

Automatically recover from failure

Scale horizontally to improve availability

Stop guessing capacity

Manage change through automation



RELIABILITY – THE BASICS

The 9's

Hard dependencies

Redundant dependencies

AVAILABILITY

Fault isolation zones

Redundant components

Microservice architecture

Recovery Oriented Computing

Distributed systems best principles

This requires that you build isolated cells of functionality within your application to prevent failure in one part taking down the entire application. You need to make sure that critical inter-cell communication is kept to a minimum, this reduces the blast radius of a cell is limited to that cell only.

DISTRIBUTED SYSTEMS BEST PRACTICES

Throttling

Retry with exponential fallback

Fail fast

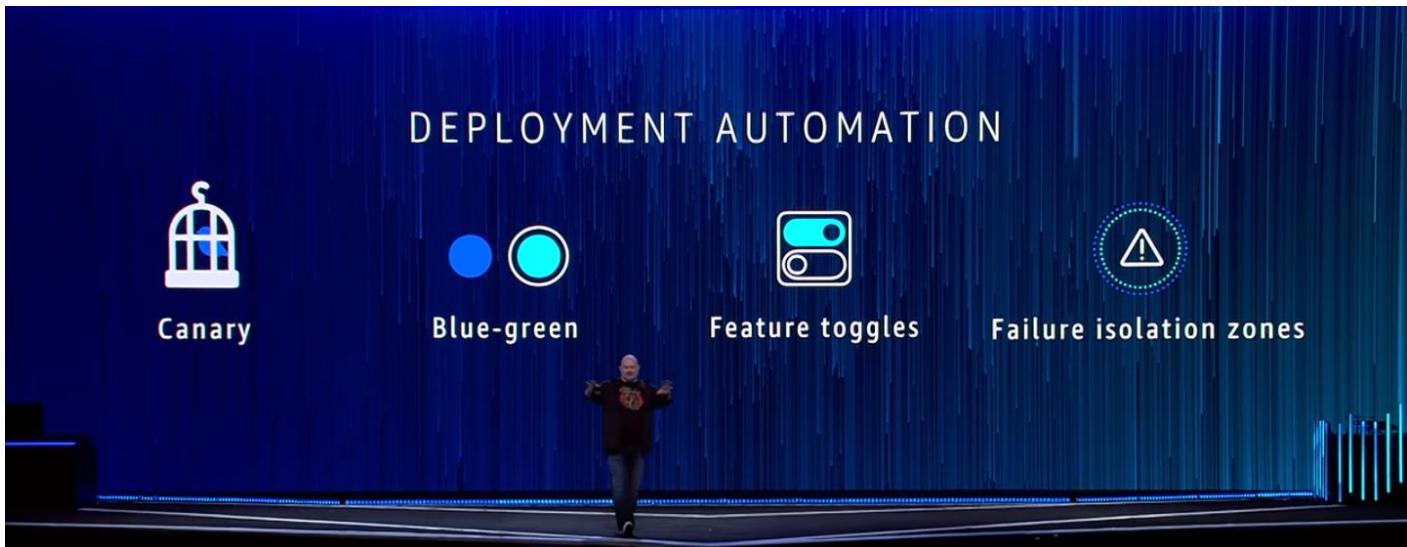
Use of idempotency tokens

Constant work

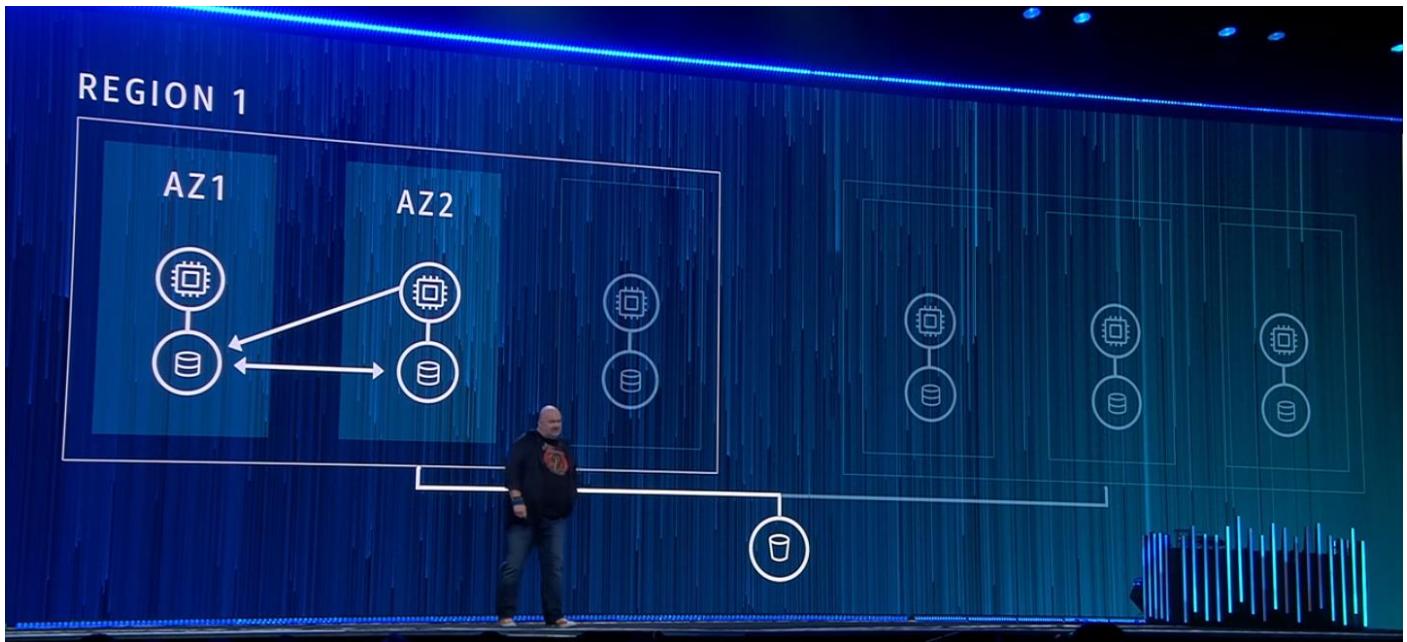
Circuit breaker

Bimodal behavior and static stability

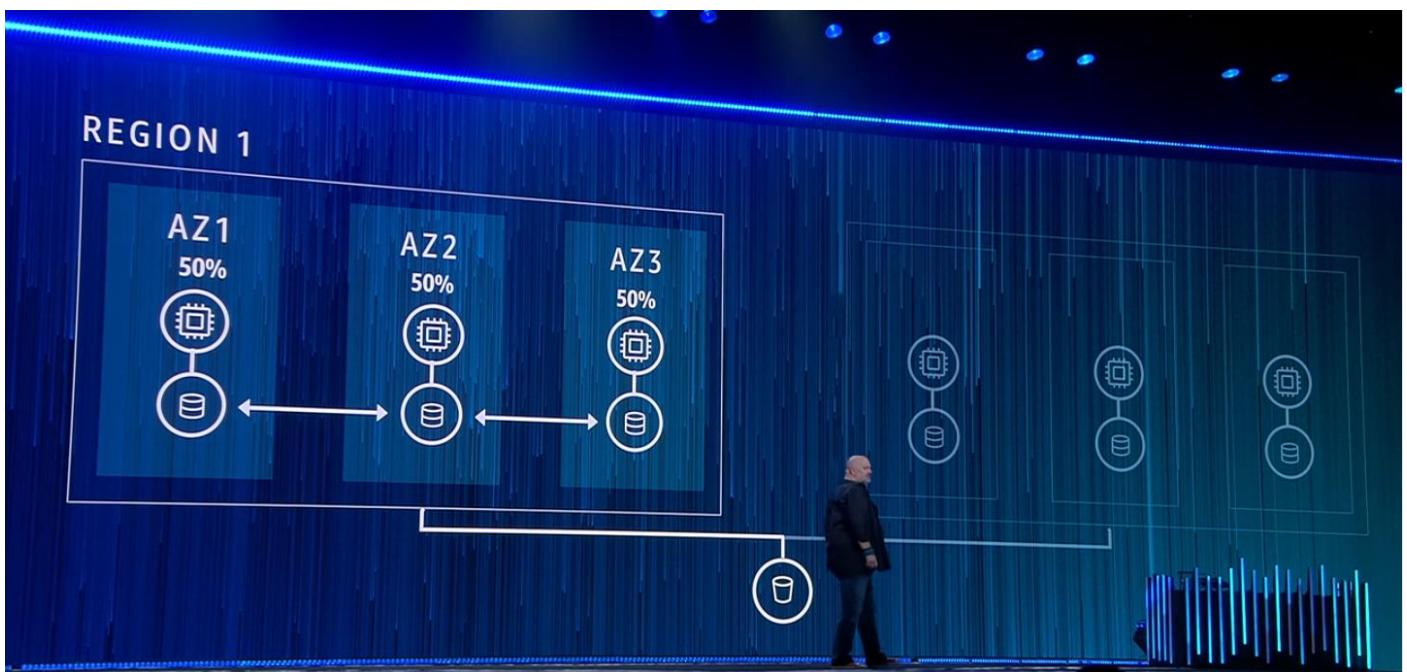
Throttle when you get overwhelmed, don't try to satisfy all calls during a load spike. Fail fast if you are going to fail so that you release resources back quickly, use circuit breakers to take bad nodes out of the loop.



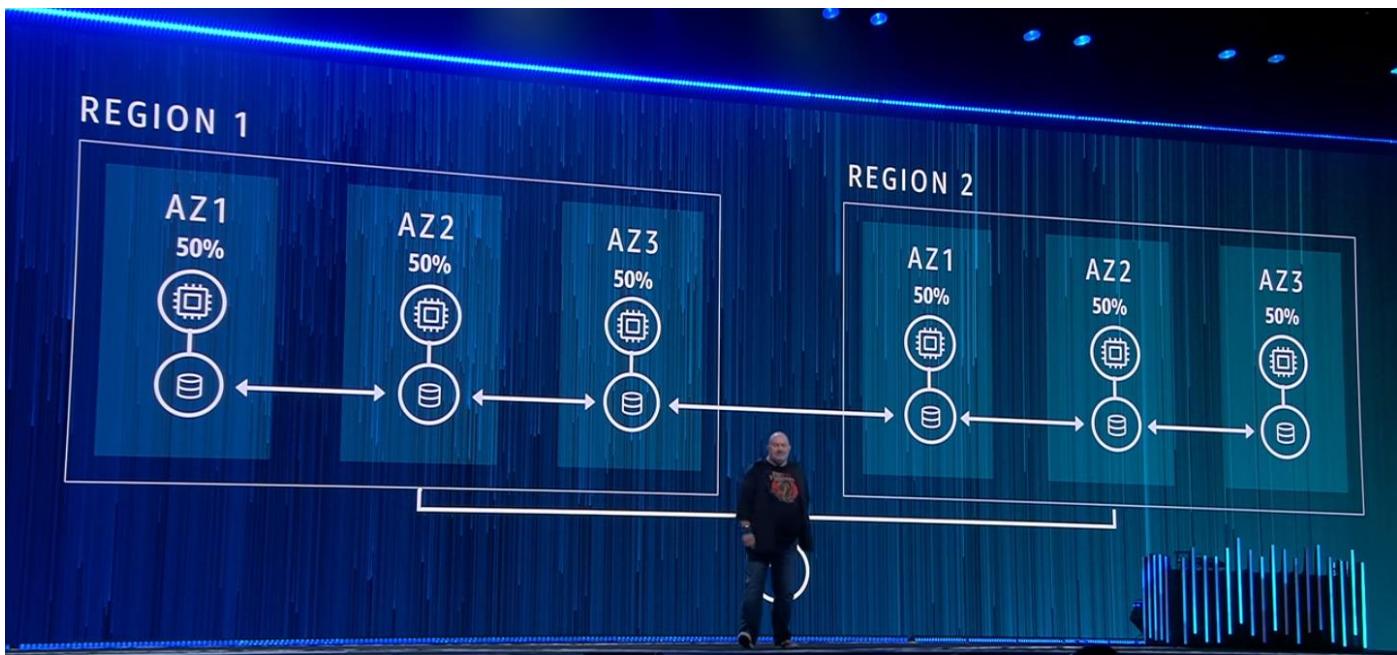
This is cheap but gives 99% availability, outages are costly and longer though



This architecture will give you a 99.5% availability when you deploy into 2 availability zones for fail overs. This limits your down time to about 30 minutes in a year.



This will give you a 99.99% availability when you deploy into 2 availability zones for fail overs. This means that if you lose an entire system in an AZ, you will still have 100% of your resources available to your users. This limits your downtime to just a few minutes in a year.



For 99.999% availability, you can deploy this active-active configuration over 2 regions and make use of route53 and use DynamoDB Global tables to give you replication of data over multiple regions. This gives you no downtime in a year availability scenario.

EXAMPLE AWS SERVICE AVAILABILITY DESIGN TARGETS	SERVICE	CONTROL PLANE	SINGLE AZ DP	MULTI AZ DP	SERVICE
	EC2	99.95	99.95	99.99	
	RDS	99.95	99.95	99.99	
	ECS	99.95			99.99
	EBS	99.95			99.999
	Kinesis				99.99
	DynamoDB				99.99
	Cloudfront	99.9			99.99
	Route 53	99.95			100
	IAM	99.9			99.995
	KMS	99.99			99.995

Check these numbers out in the reliability tables within the Well Architected Framework document

TEST, TEST, TEST

Break your systems and see how they respond to failures, there is a lot to be learnt about your system architecture



Suppose a section of your screen isn't loading its data, this should result in us not showing you that section anymore and instead show the other things available



Whoops, something went wrong...

Netflix Streaming Error

We're having trouble playing this title right now. Please try again later or select a different title.

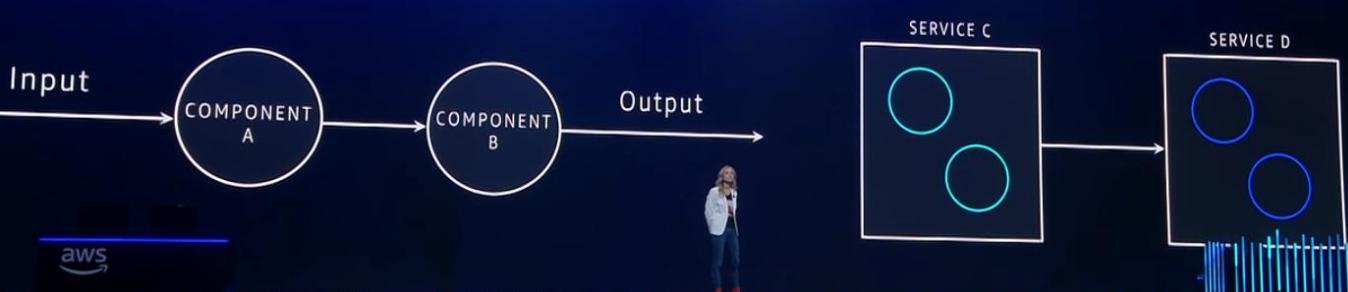


The service failing should not result in this type of screen ever! We use chaos experiments to guarantee this never happens.

UNIT TEST

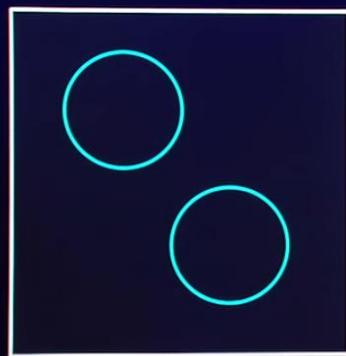


INTEGRATION TEST

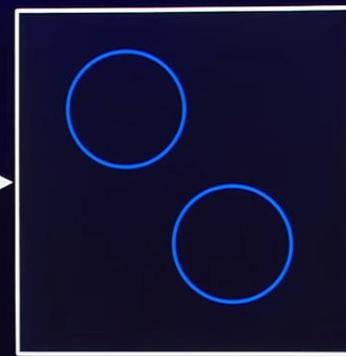


CHAOS EXPERIMENTS

SERVICE C



SERVICE D



In **Chaos experiments**, there is an option to add a failure or add a time latency to the response going back to the calling service between calls. Here we have unknowns and are resilient to these failures

A large, vibrant background graphic featuring abstract, swirling lines in various colors like yellow, green, red, and blue. Overlaid on this graphic is the text "CHAOS ENGINEERING" in a large, white, sans-serif font.



This is about experimenting on a distributed system in order to build confidence in the systems capability to withstand turbulence in production. Chaos engineering is not replacing other types of tests, it is for guaranteeing resilience

NETFLIX



O'REILLY® Chaos Engineering

Building Confidence in System Behavior
through Experiments

Casey Rosenthal, Lori Kochstein,
Aaron Blohowiak, Nic Jones
& Ali Basiri





This is also known as Chaos monkeys, bringing down a single instance of some service

TARGETED CHAOS

aws

FORCES OF CHAOS

Graceful
restarts +
degradation

Targeted
chaos

aws

CAN WE CAUSE A CASCADING FAILURE?

aws

FORCES OF CHAOS

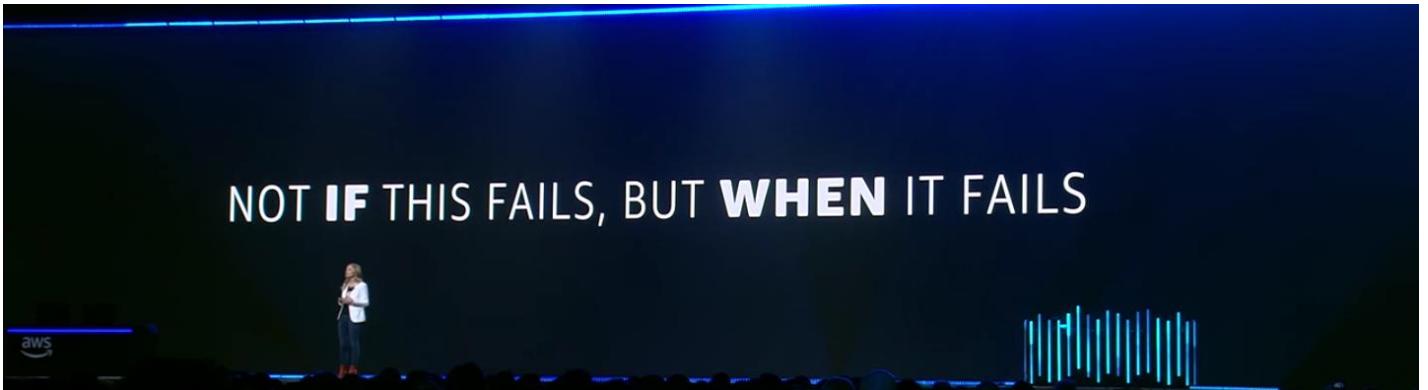
Graceful
restarts +
degradation

Targeted
chaos

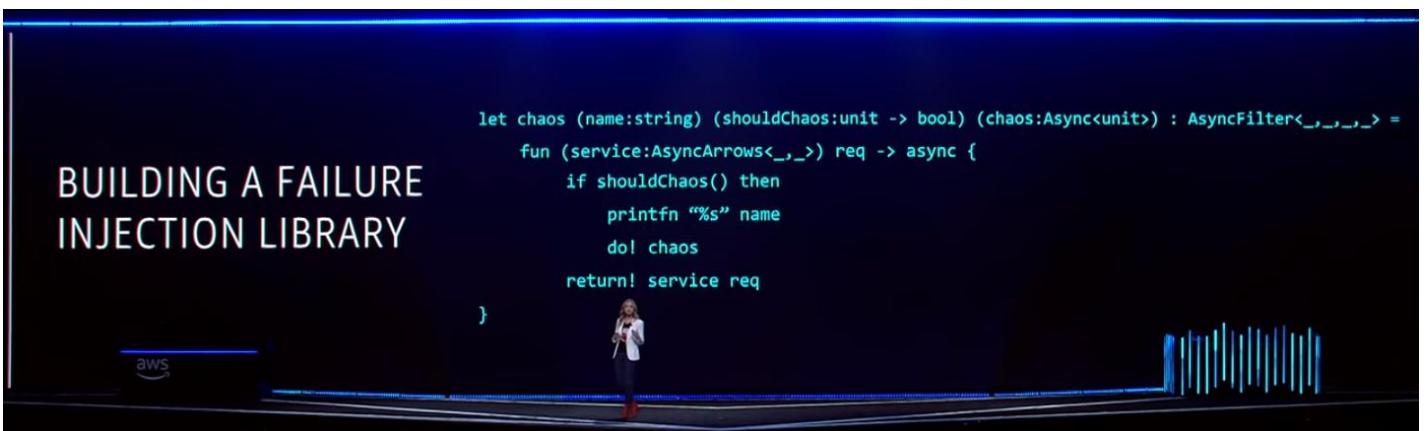
Cascading
failure

aws

This is when failure in some part of the system triggers failures in other parts, then in other systems. We tried to cause cascading failures and watch what happens on some monitoring screen. We are doing this in QA and NOT production!



Netflix does failure injection testing where service owners can add latency or failures to the code for resiliency testing if predefined criteria conditions are met



SAFETY & MONITORING



With the failure injection testing framework, Netflix has built an automated chaos monitoring platform called **Chap** to determine how much services traffic they are impacting at any moment in time

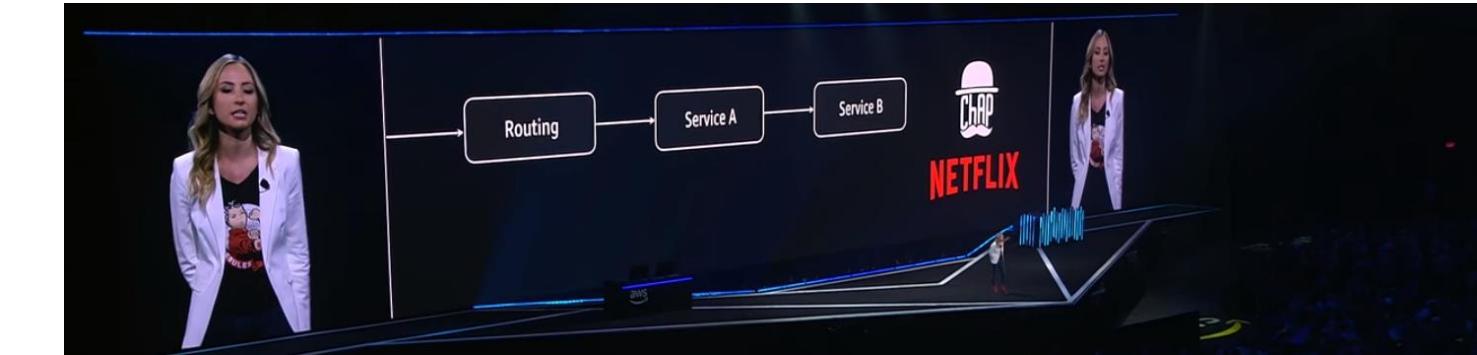
FORCES OF CHAOS

Cascading failure

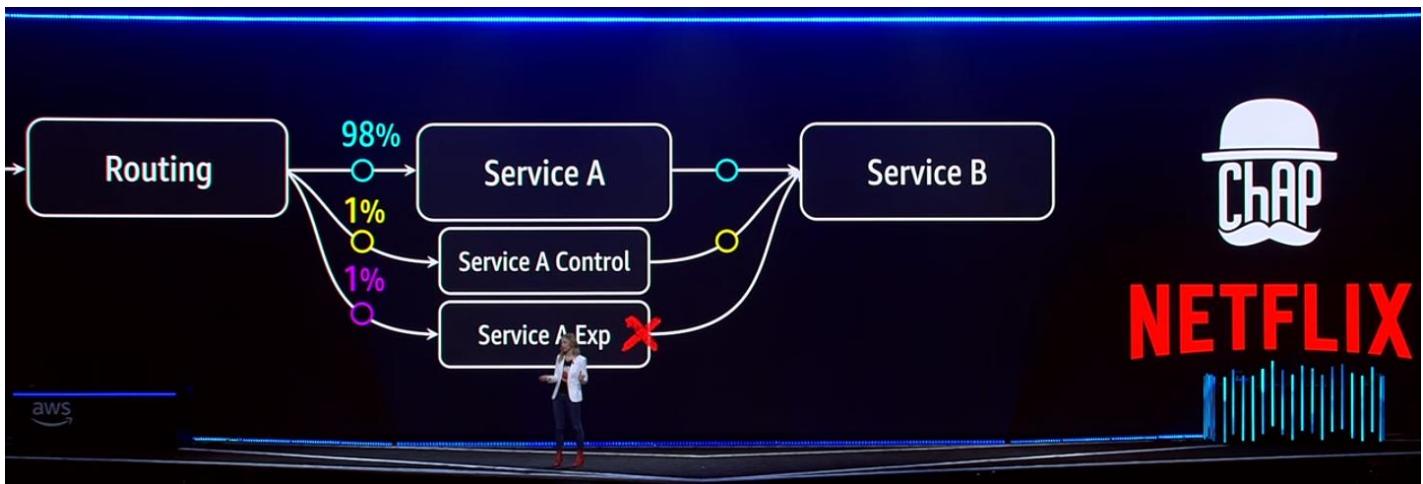
Failure injection

Chaos automation platform

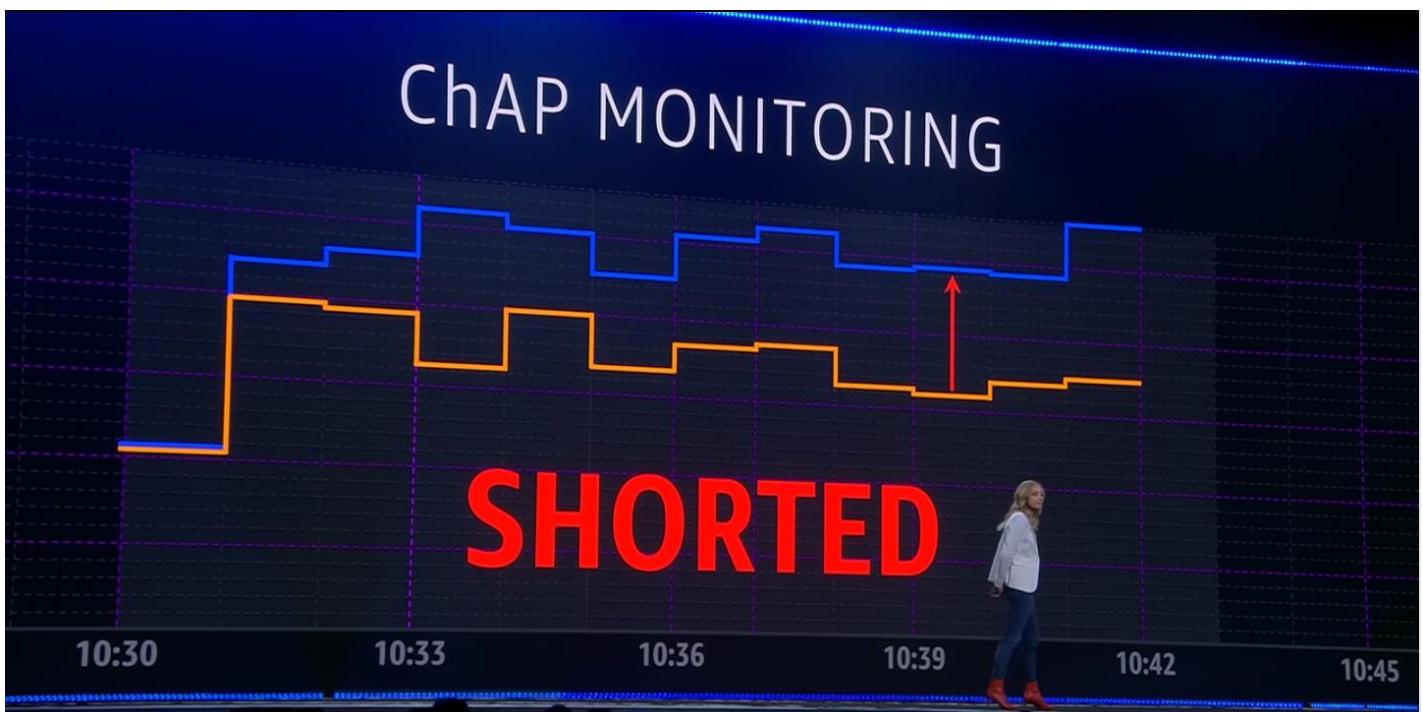
ed



The **key business metric at Netflix is whether or not the user can press 'Play'** and nothing starts. This metric is calculated during chaos experiments



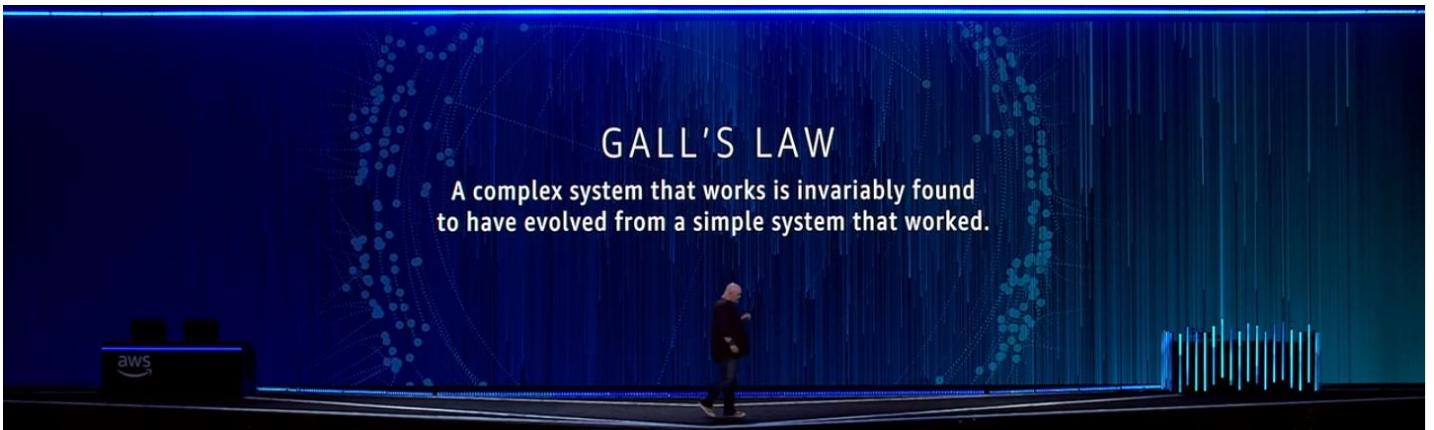
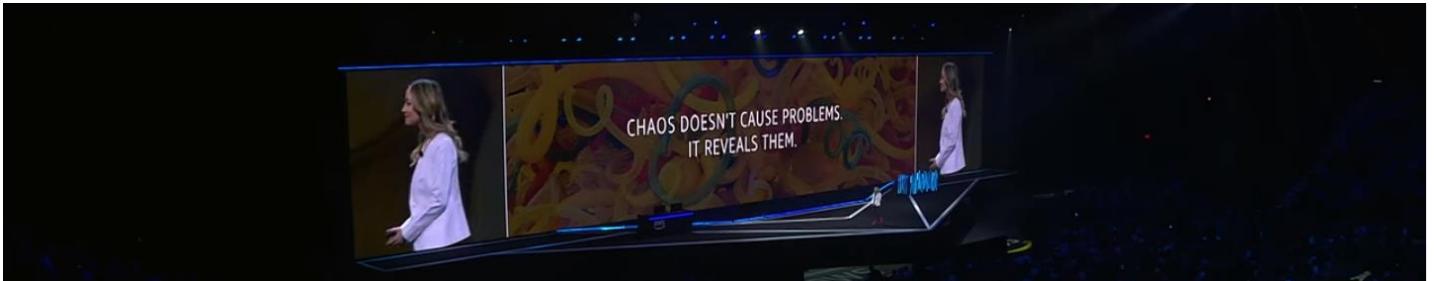
We create 2 clusters for control and experiments and watch the KPI



Too much deviation means that we stop the experiments ongoing and engineers can start to debug the issue



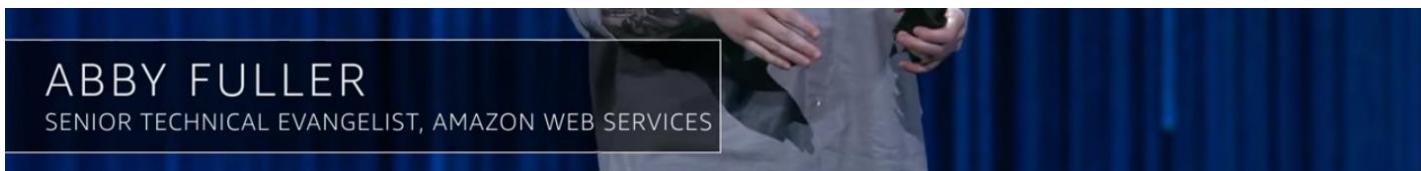
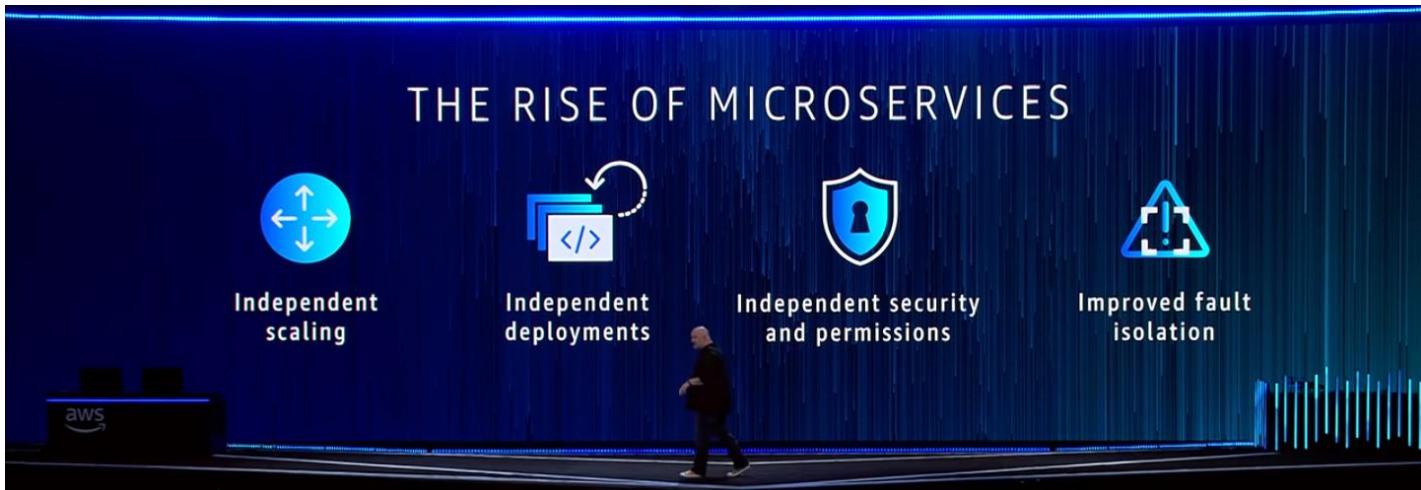
We now have an automated chaos experiments generating algorithm, this determines what services needs to be tested in certain scenarios.



You need to start off by building smaller, simpler systems and then aggregate them into larger a complex system



The more managed services you use, the less moving parts you have regarding the control and data planes



CONTAINERS ON AWS



Amazon ECS



AWS Fargate for ECS



Amazon EKS



COMING SOON
AWS Fargate for EKS



THE **POWER** IS
IN THE **CHOICES**



BUT HOW CAN YOU **GET STARTED?**



LET'S **REWIND**



On AWS since 2014 with ECS







CAPITAL ONE RUNS SERVICES ON
ECS AND **DOCKER** FOR A FASTER
PACE OF DEVELOPMENT.



MARKETING DATA STARTUP
MAJORITY OF SERVICES
RUNNING ON ECS



25,000 TASKS
382 SERVICES

Average 100s of containers
per host with binpacking



ECS





ECS IS NOT JUST FOR
EVENT PROCESSING



RUNNING ETL JOBS ON ECS

aws

**ECS ISN'T THE ONLY WAY
TO RUN CONTAINERS IN PRODUCTION**

aws



MORE CUSTOMERS RUNNING
KUBERNETES ON AWS THAN
ANYWHERE ELSE

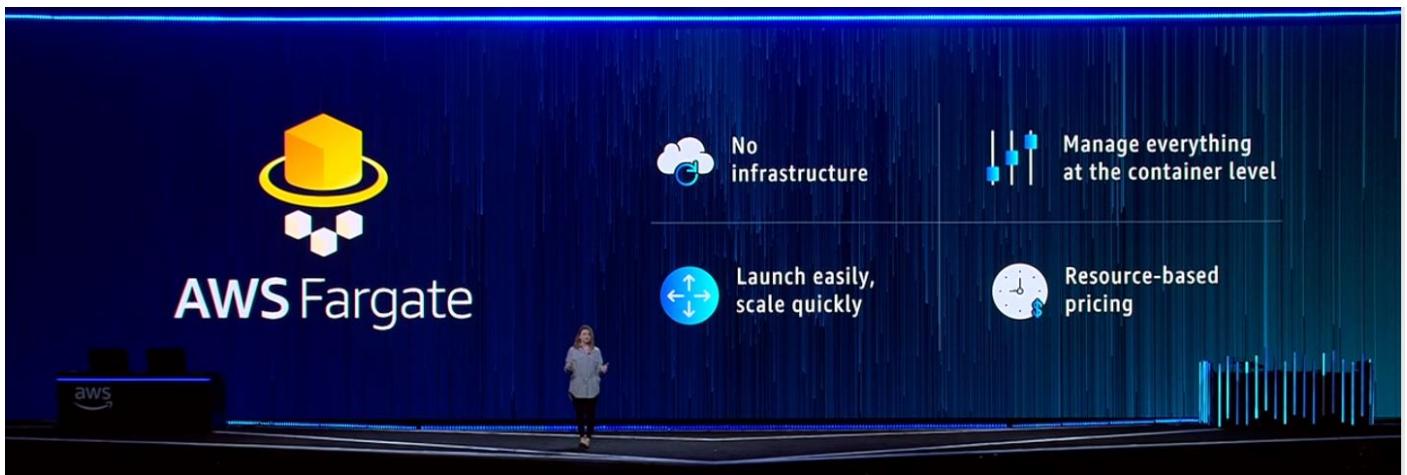
aws



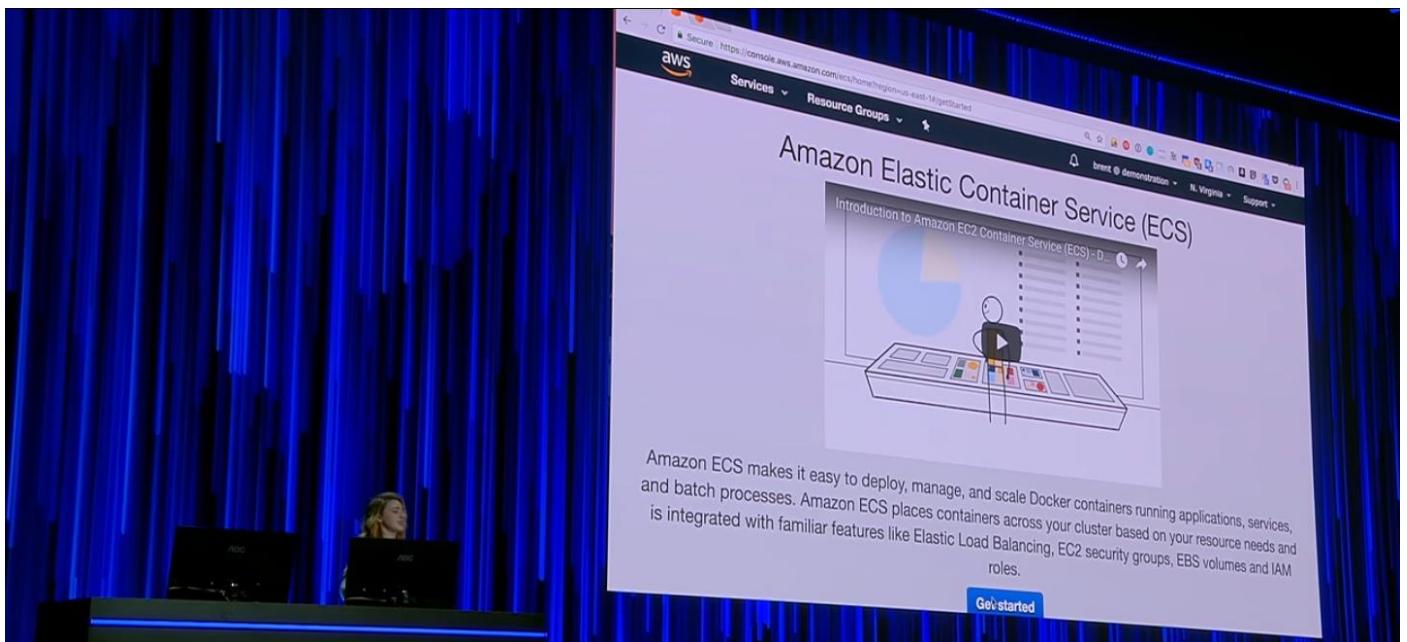


"RUN KUBERNETES FOR ME"

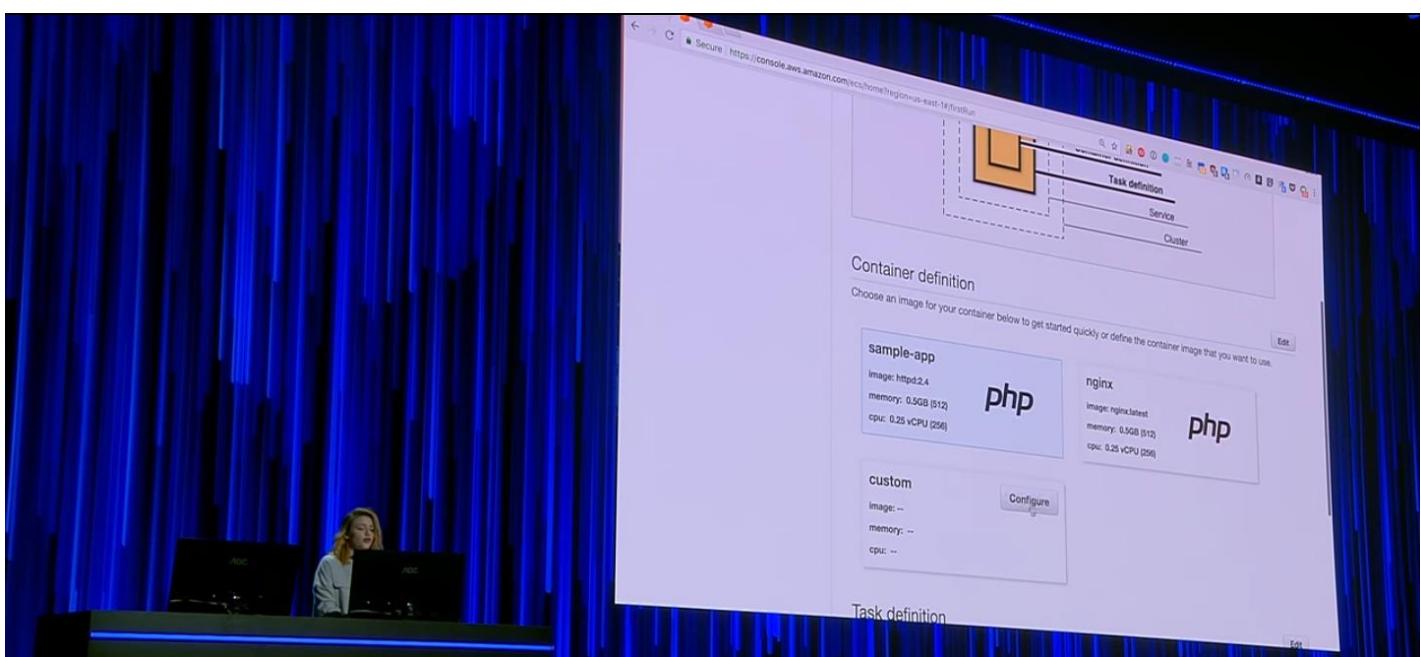
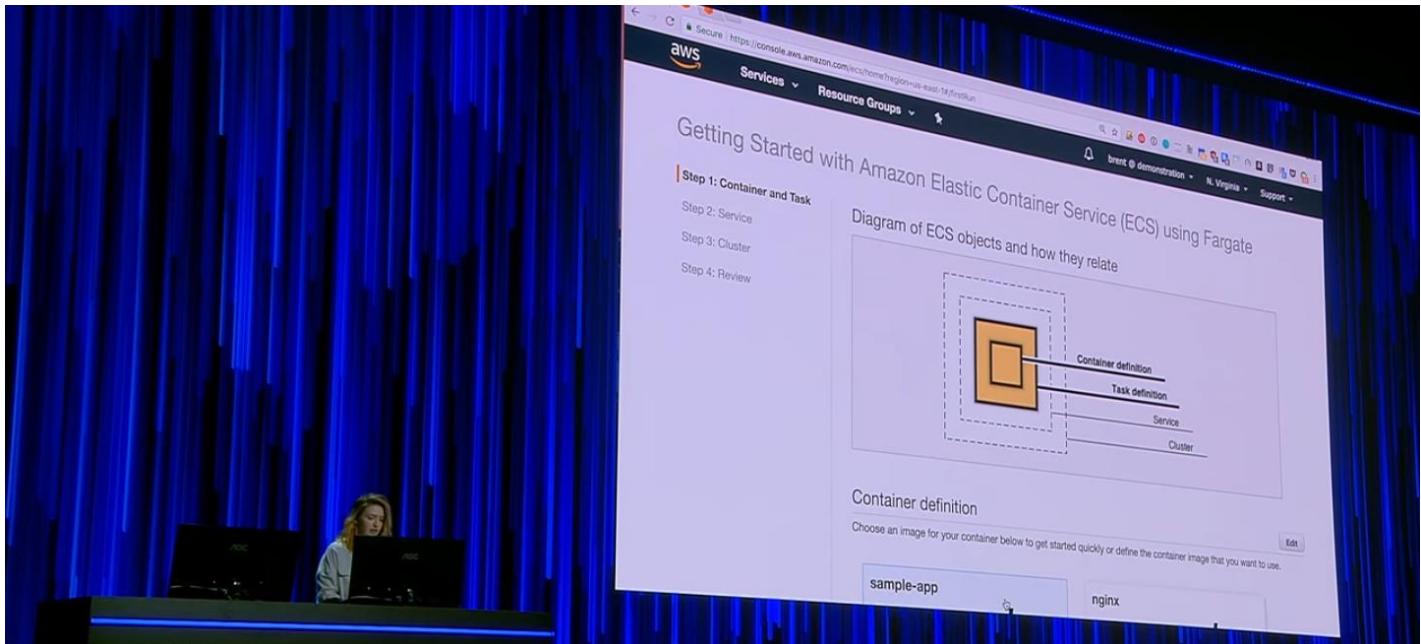
A large, stylized quote "RUN KUBERNETES FOR ME" is displayed prominently in white capital letters against a dark blue background. A woman stands on the stage to the right of the quote, holding a microphone and gesturing towards it. The monzo logo is visible in the bottom left corner.



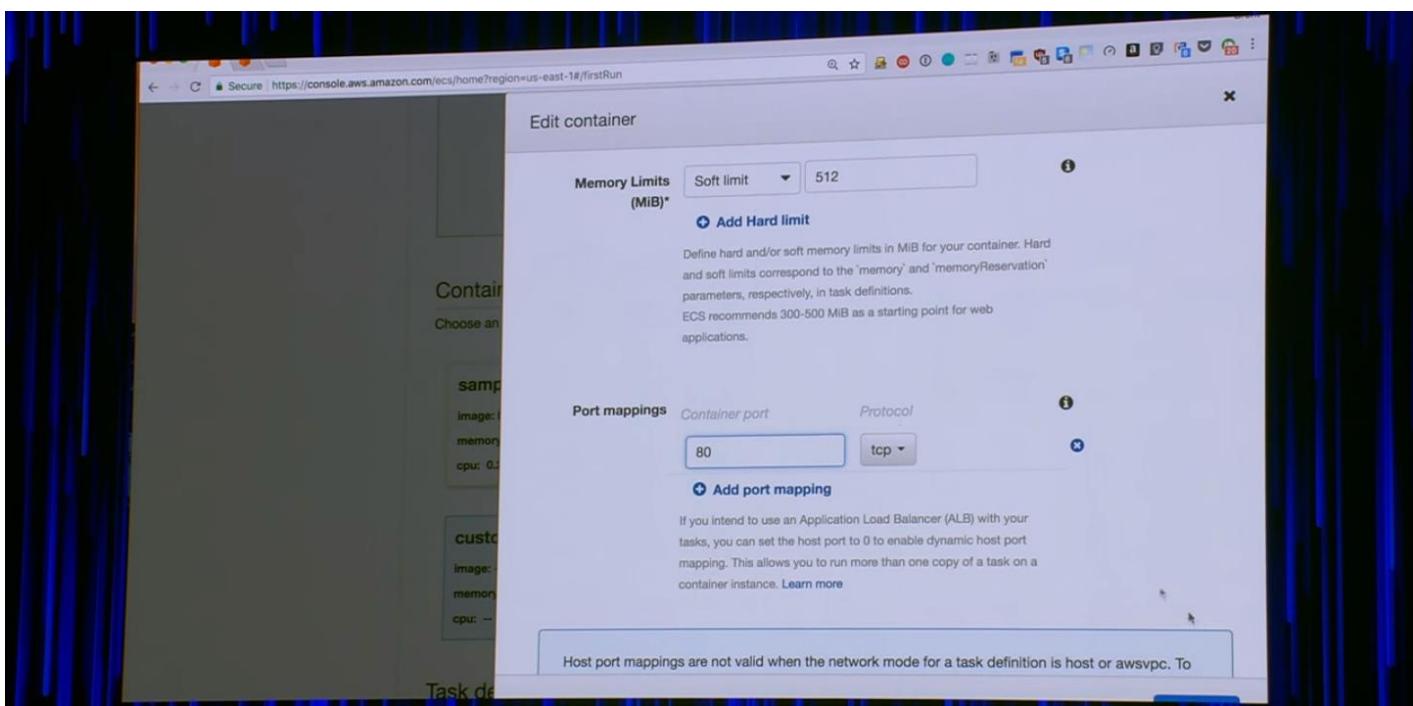
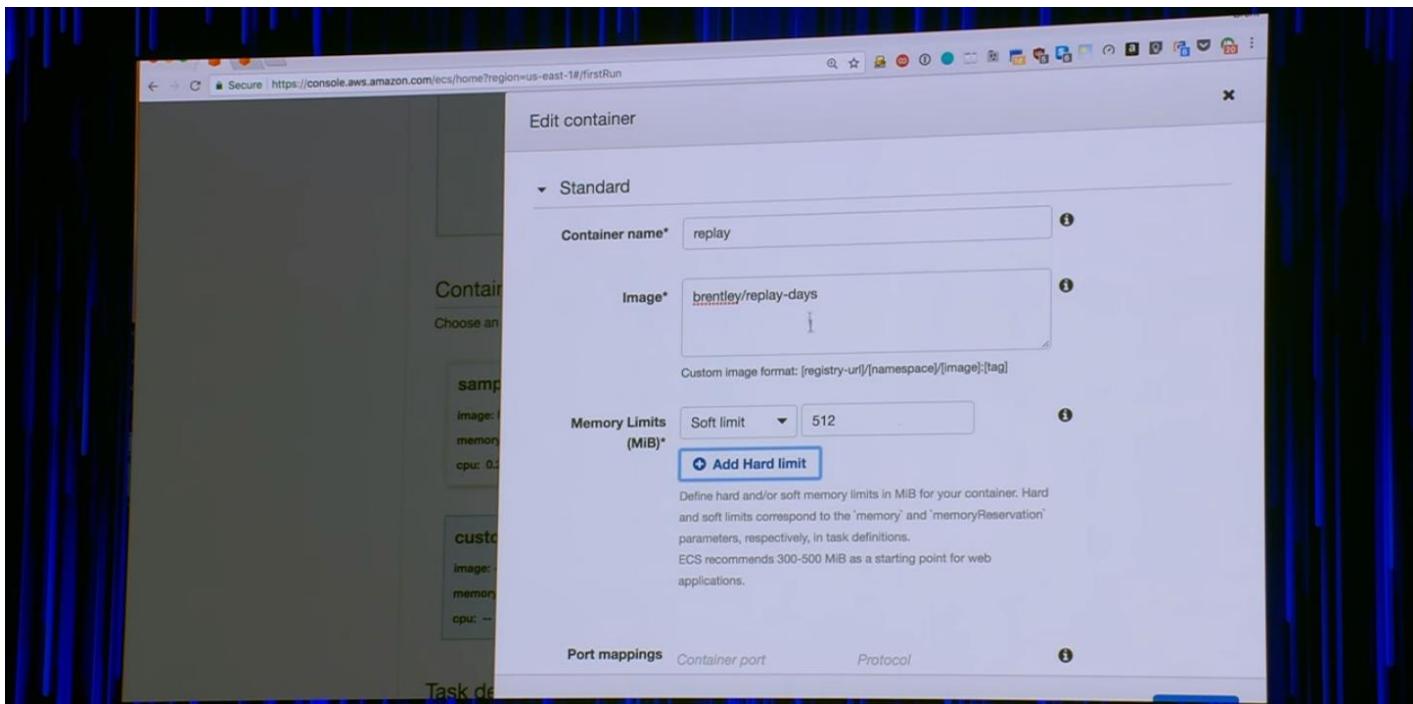
AWS Fargate is not a service but is an underlying technology that you can use to power your container workloads. Just take your task definition or Pod, drop it into Fargate, add some resources, and away you go.



Let's see a demo of how to deploy an app and see Fargate in action.



We start with our task, within the Fargate experience



The screenshot shows the AWS CloudWatch Metrics console. A task definition named "replay" is selected. The task definition details are as follows:

- Task definition name: first-run-task-definition
- Network mode: awsvpc
- Task execution role: Create new
- Compatibilities: FARGATE
- Task memory: 0.5GB (512)
- Task CPU: 0.25 vCPU (256)

An "Edit" button is located in the top right corner of the task definition section.

The screenshot shows the "Configure task definition: replay" dialog box. The task definition details are:

- Task definition name*: replay
- Network mode*: awsvpc
- Task execution role: ecsTaskExecutionRole
- Compatibilities*: FARGATE

The "Task size" section is visible below, showing a dropdown for "Task memory*" set to 0.5GB (512).

The screenshot shows the AWS CloudWatch Metrics service configuration page. At the top, there's a search bar and a navigation menu. Below the header, a card displays the task definition details:

replay Configure

image: brentley/replay-days
memory: 0.5GB (512)
cpu:

Task definition Edit

A task definition can be thought of as a blueprint for your application that describes one or more containers. Some attributes are configured at the task level but the majority of attributes are configured per container.

Task definition name: replay Info

Network mode: awsvpc Info

Task execution role: Create new Info

Compatibilities: FARGATE Info

Task memory: 0.5GB (512) Info

Task CPU: 0.25 vCPU (256) Info

The screenshot shows the 'Define your service' step of the AWS ECS service creation wizard. A woman is visible in the background, sitting at a desk with two monitors. The wizard interface includes a back button and a 'Next Step' button.

Define your service Edit

A service allows you to run and maintain a specified number (the "desired count") of instances of a task definition simultaneously in an ECS cluster.

Service name: replay-service Info

Number of desired tasks: 1 Info

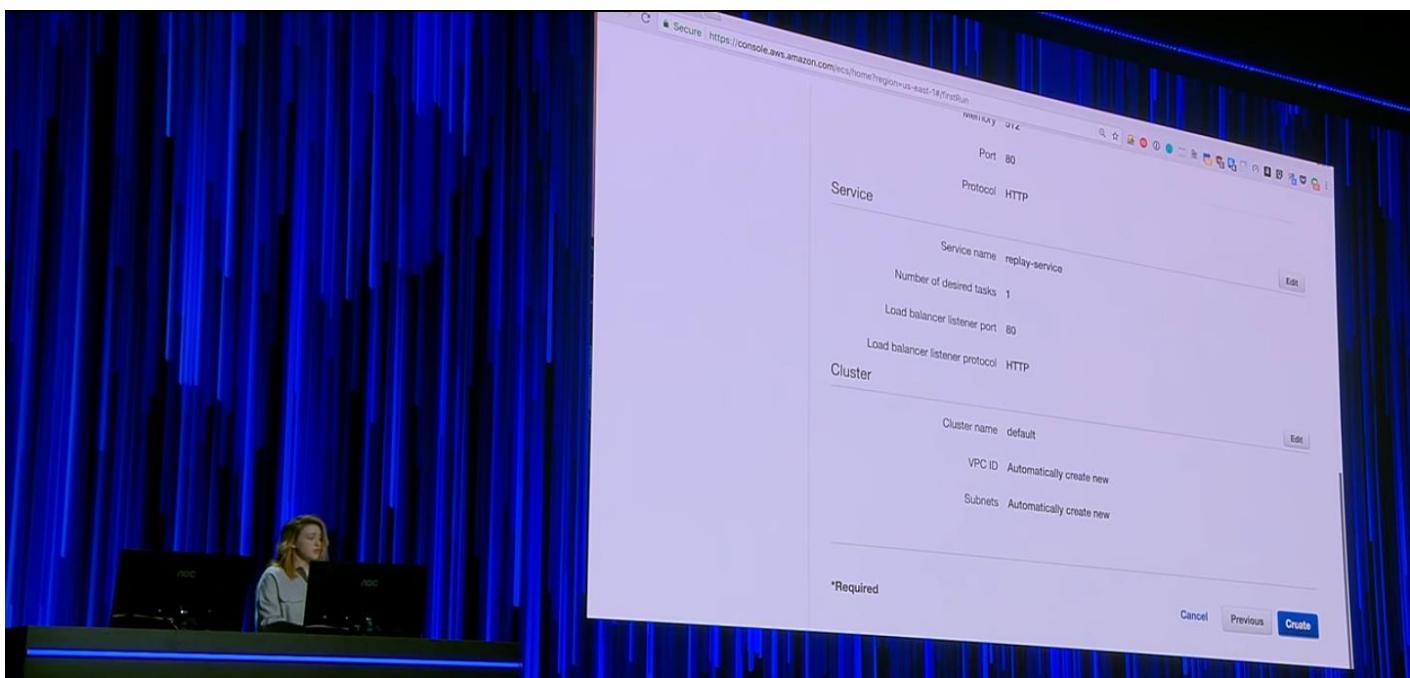
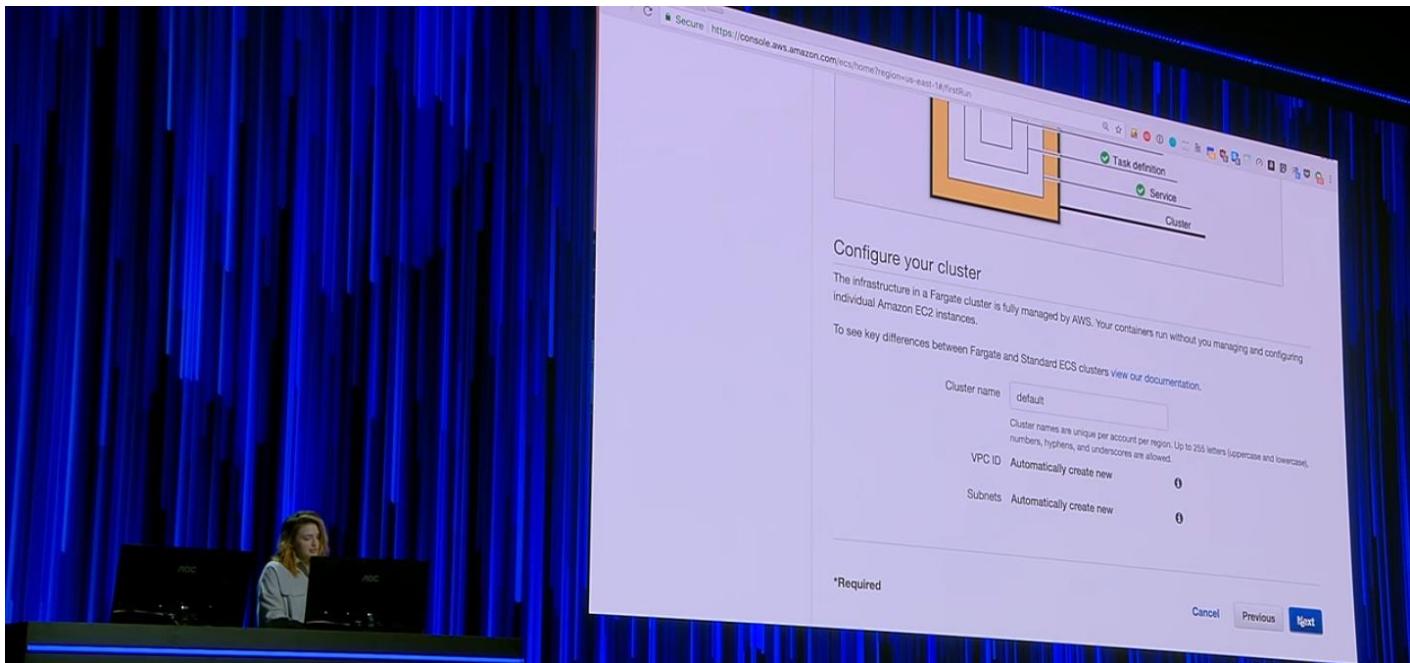
Security group: Automatically create new Info
Two security groups will be created to secure your service. An ALB security group will be created that allows all traffic on the ALB port and an ECS security group will be created ONLY from the ALB security group. You can further configure security groups and network access outside of this wizard.

Load balancer type: Application Load Balancer (ALB) None

Load balancer listener port: 80

Load balancer listener protocol: HTTP

*Required Cancel Previous Next



Secure | https://console.aws.amazon.com/ecs/home?region=us-east-1#/firstRun

Launch Status

Creating resources for your service, this may take up to 10 minutes. When we're complete, you can view your service.

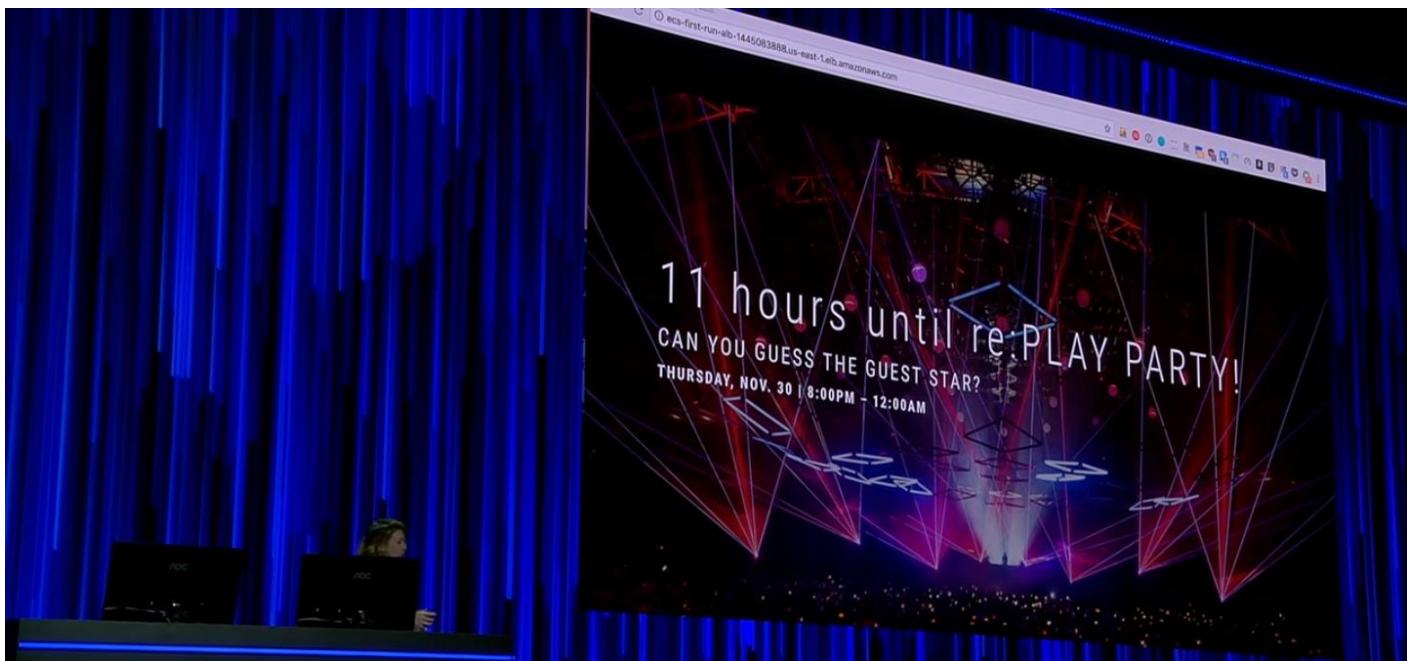
Back [View service](#) Enabled after service creation completes successfully

Additional features that you can add to your service after creation

Scale based on metrics
You can configure scaling rules based on CloudWatch metrics

Preparing service: 2 of 10 complete

ECS resource creation	Status
Cluster default	complete ✓
Task definition replay:3	complete ✓
Service	pending ⏱
Additional AWS service integrations	
Log group Log group [/ecs/replay] already exists	complete ✓
CloudFormation stack	pending ⏱
VPC	pending ⏱
Subnet 1	pending ⏱
Subnet 2	pending ⏱
Security group	pending ⏱
Load balancer	pending ⏱



The app is now deployed.

IT'S NOT
ABOUT HOW

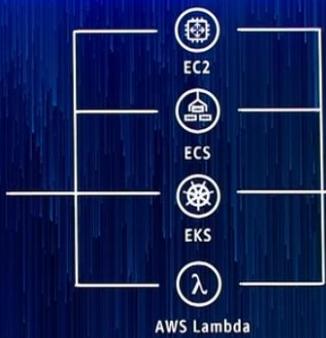


IT'S ABOUT
HOW WELL

aws

```
var http = require('http');
http.createServer(function (req, res) {
  res.writeHead(200, {'Content-
  Type': 'text/plain'});
  res.end('Hello World!');
}).listen(8080);
```

BUILD



AWS Fargate lets you focus on the workload only

"RUN MY APPLICATION."

aws

We will be able to build with a couple of basic primitive – the container

SECURE. SCALABLE. RELIABLE.



WORKLOAD FIRST, FULL STOP.