# AWS re:Invent

DECEMBER 2 – 6, 2024 | LAS VEGAS, NV

ARC311

# Scaling Prime Video for peak NFL streaming on AWS

**Elliott Nash**

(He/Him)
Head of Infrastructure
Scaling
Prime Video

**Ralph Chaker**

(He/Him)
Principal Product Manager
Prime Video

**Tulip Gupta**

(She/Her)
Sr. Solutions Architect, AWS
Strategic Accounts
AWS

# Agenda

- Introduction to NFL Thursday Night Football on Prime Video

- Prime Video multi-Region scaling strategies

- Prime Video elastic scaling strategies

- AWS benefits – AWS Well-Architected Framework

aws

120+

**prime video**

### Business opportunity
Cost-effectively scale and delight millions of subscribers around the world during live sports

### Challenge
Scale out infrastructure to meet the demands of NFL games viewership

### Solution
Adopting autoscaling and multi-Region architecture

# FAQs: Thursday Night Football on Prime Video

**1** September 7, 2017, the first NFL game was co-exclusively streamed on Prime Video

**2** In 2022, exclusively streamed on Prime Video

**3** In 2023, 17 million fans watched TNF on Prime Video

**4** 2024 Season 3 average was 18% higher than 2023

**5** 2024 biggest NFL event in Prime Video history

A multi-Region
scaling story

# We are global

**240+**
Countries and territories

**30+**
Languages

**200M+**
Prime members worldwide

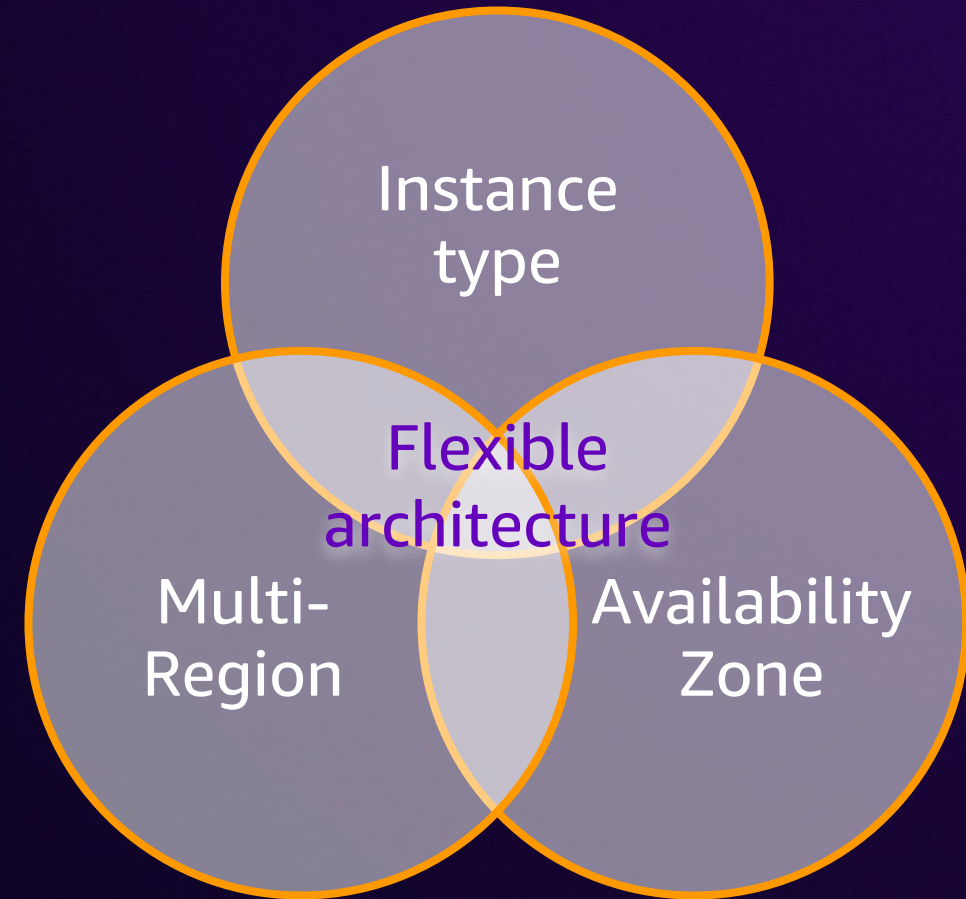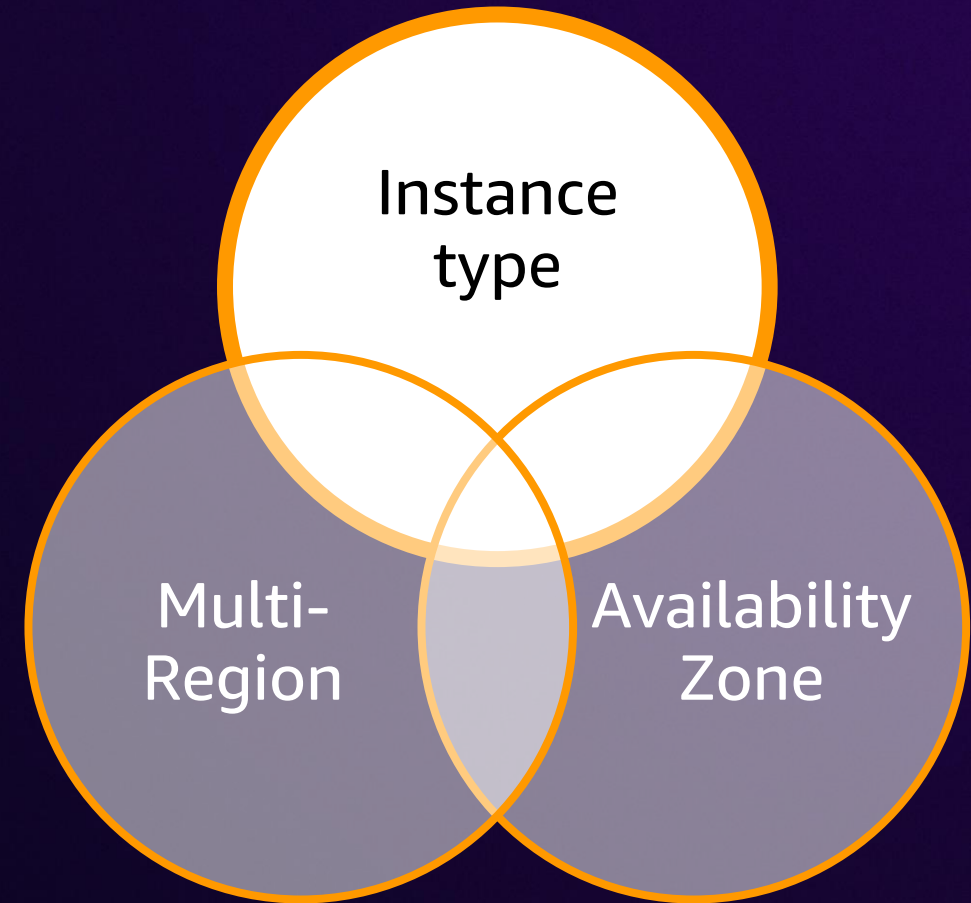| | | |
|---|---|---|
| Chinese | Indonesian | Russian |
| Danish | Norwegian | Spanish |
| Dutch | Polish | Swedish |
| English | Portuguese | Tamil |
| French | Italian | Telugu |
| German | Japanese | Thai |
| Hindi | Korean | Turkish |

prime video

# Designing for high availability

- Availability requirements:
  - Better than broadcast
  - Live content is perishable

- Flexibility across three layers

Instance type

Flexible architecture

Multi-Region

Availability Zone

# Designing for high availability

- Availability requirements:
  - Better than broadcast
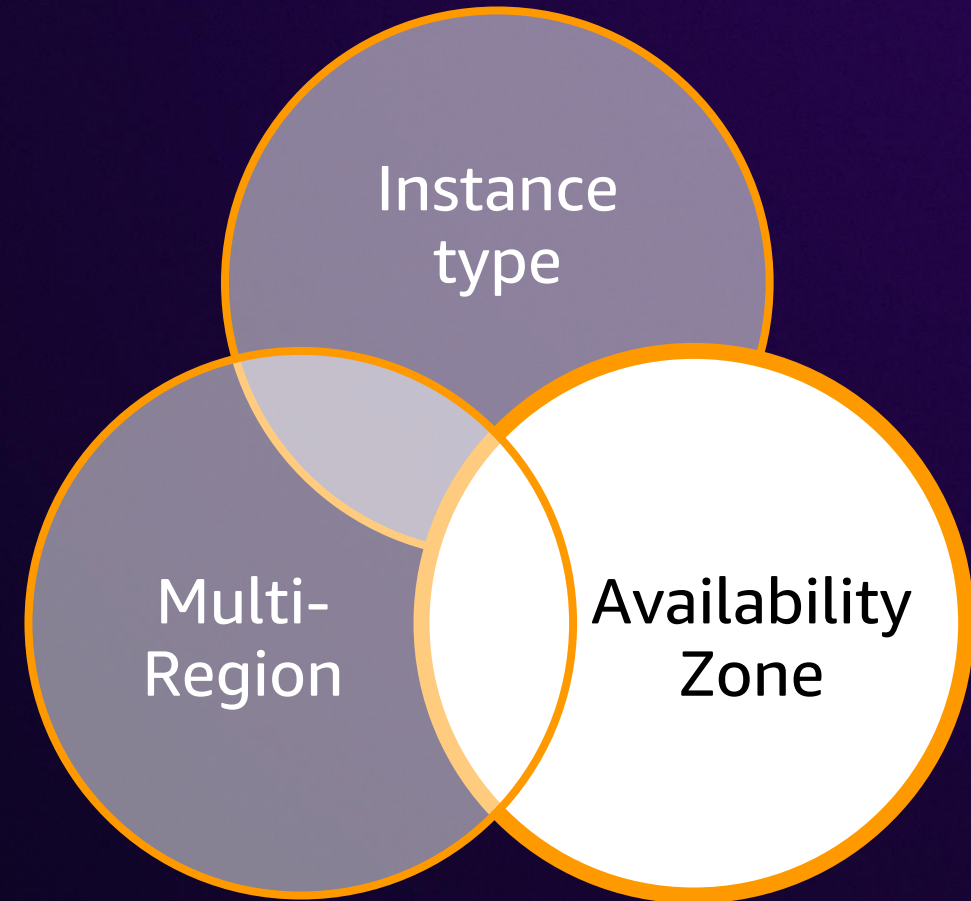  - Live content is perishable

- Flexibility across three layers

Instance type

Multi-Region

Availability Zone

# Designing for high availability

- Availability requirements:
  - Better than broadcast
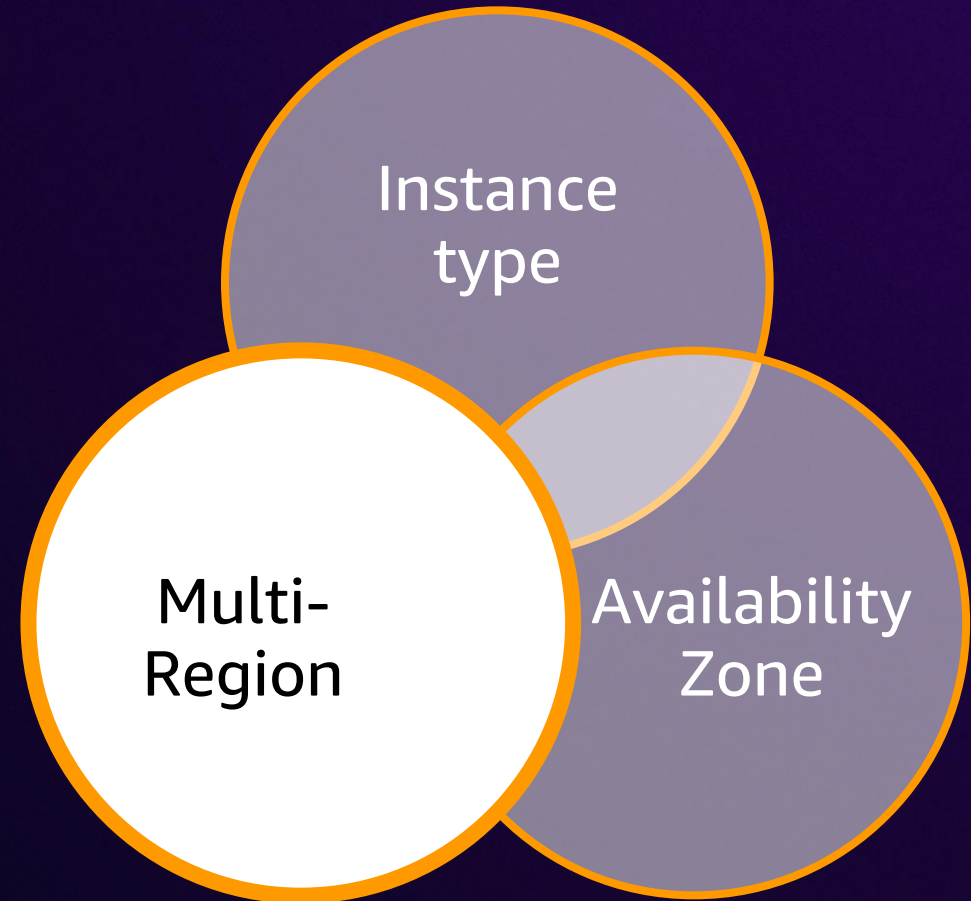  - Live content is perishable

- Flexibility across three layers



Instance type

Multi-Region

Availability Zone

# Designing for high availability

- Availability requirements:
  - Better than broadcast
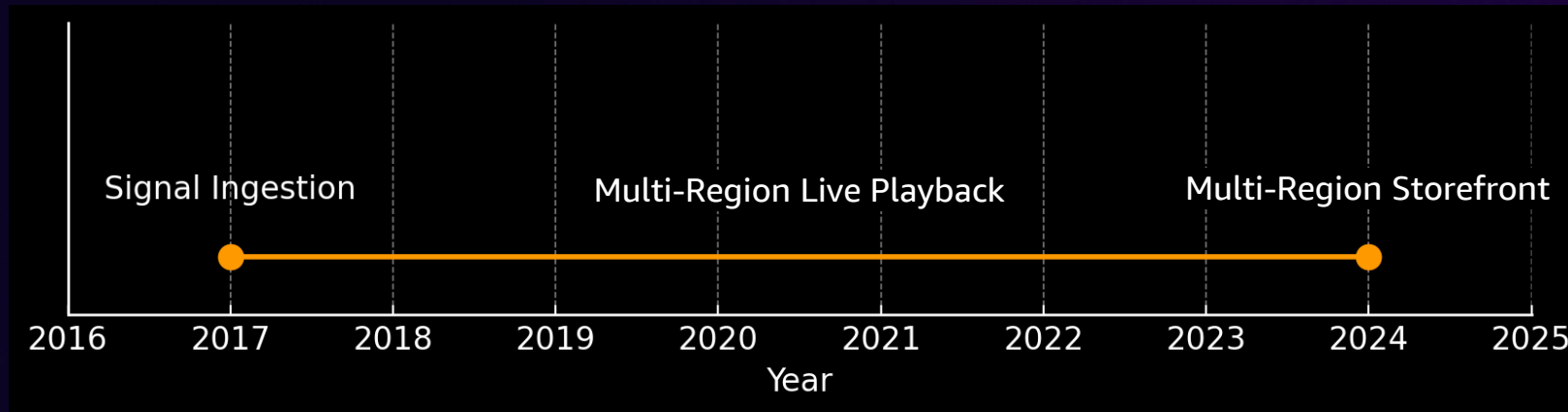  - Live content is perishable

- Flexibility across three layers

Instance type
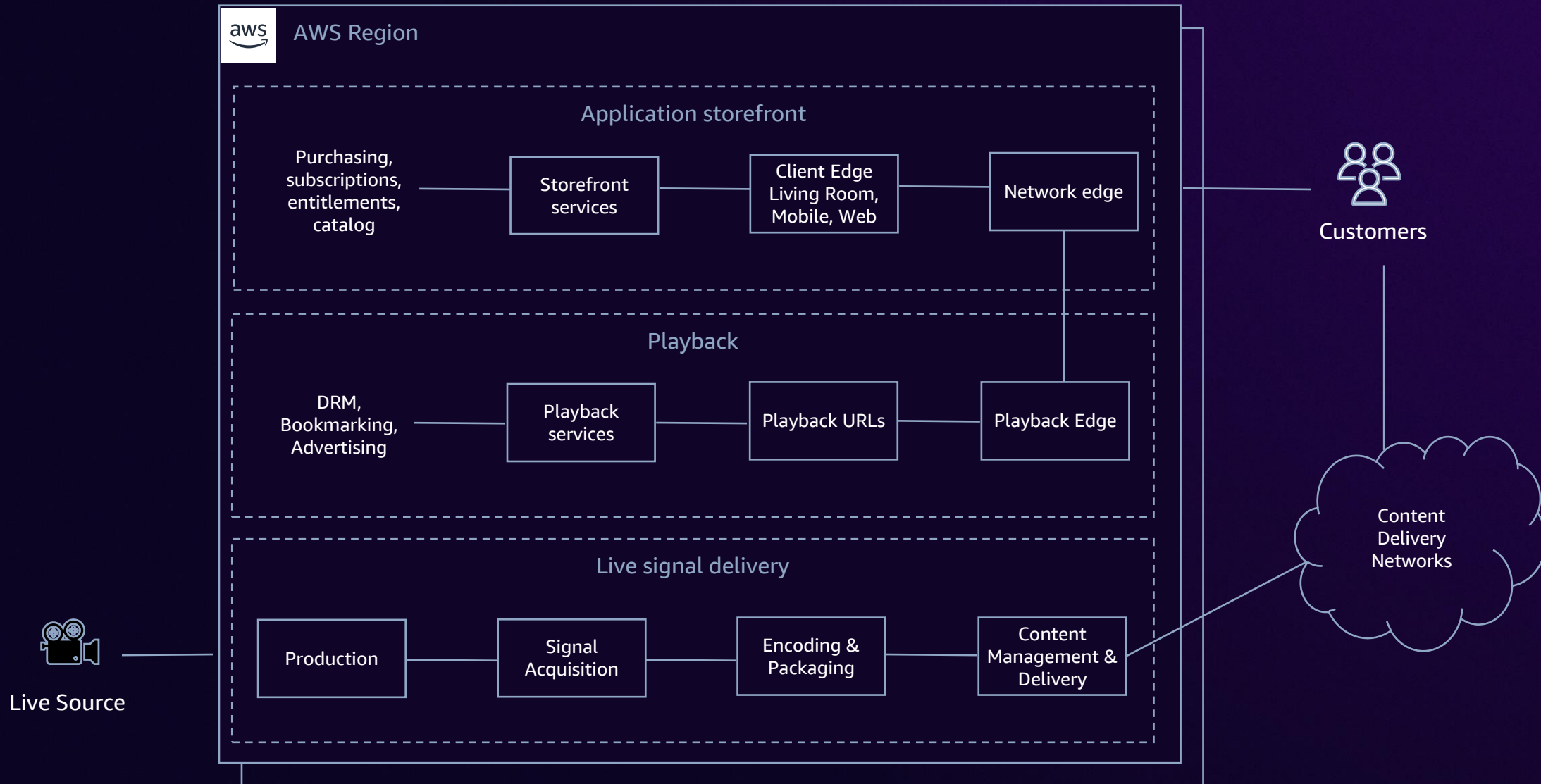
Multi-Region

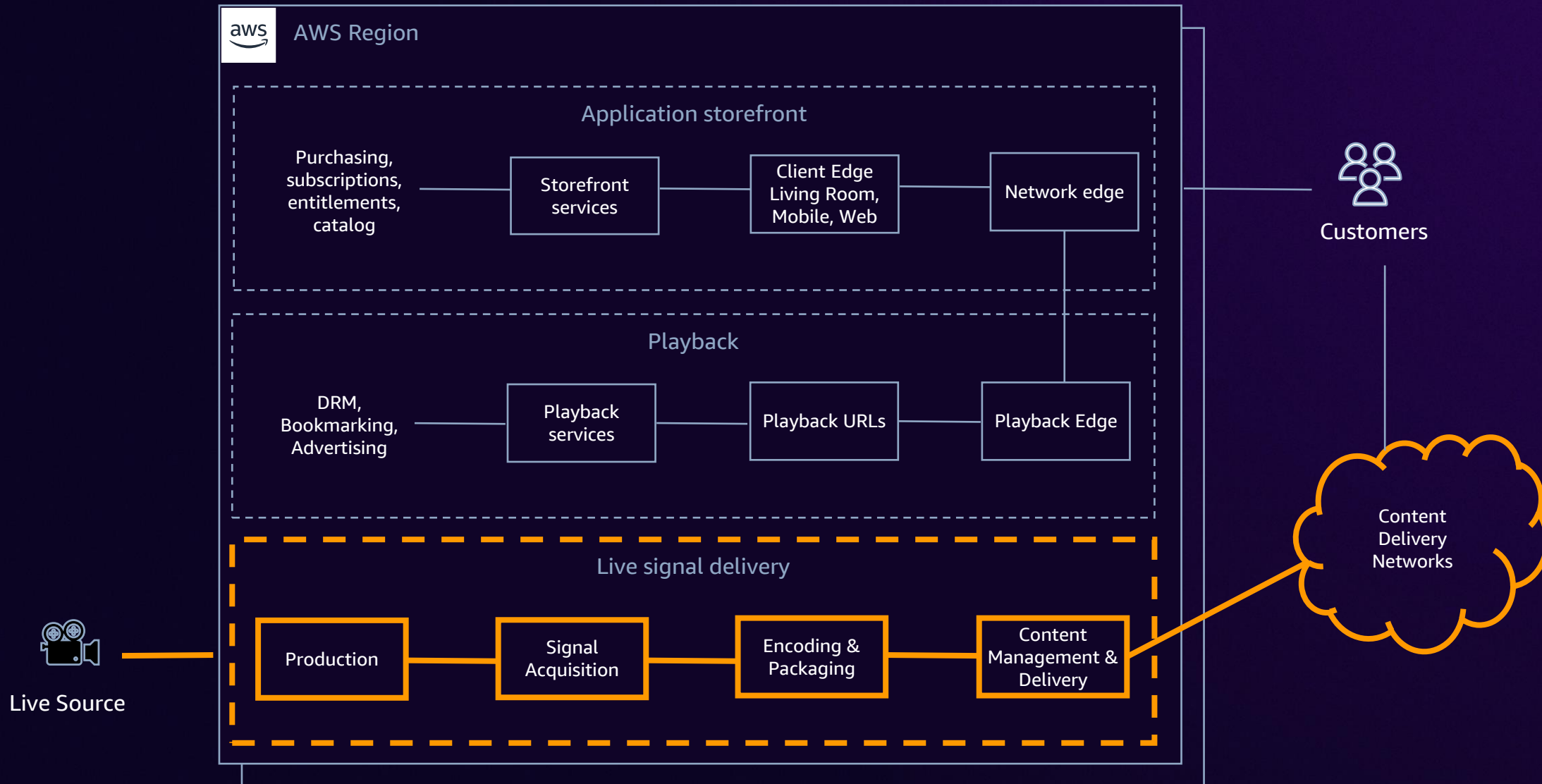Availability Zone

# Architecture evolution

- 2017: Launch TNF on Prime Video (non-exclusive)
- 2022: First TNF exclusive season
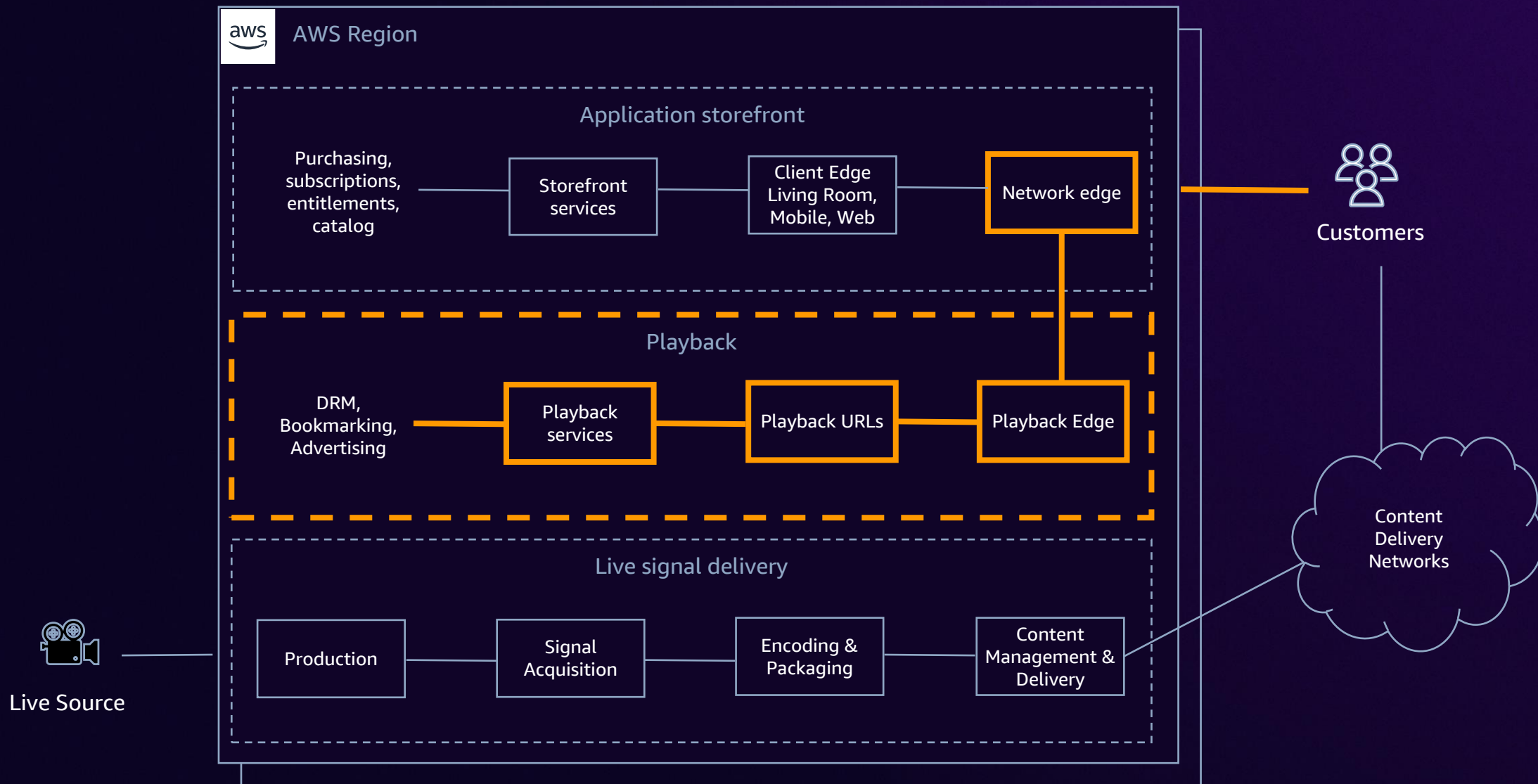- 2024: All critical customer journeys support multi-Region

# 10,000 ft architecture



**AWS Region**

### Application storefront

Purchasing, subscriptions, entitlements, catalog → Storefront services → Client Edge Living Room, Mobile, Web → Network edge

### Playback

DRM, Bookmarking, Advertising → Playback services → Playback URLs → Playback Edge

### Live signal delivery

Production → Signal Acquisition → Encoding & Packaging → Content Management & Delivery

Live Source

Customers

Content Delivery Networks

# 10,000 ft architecture

**AWS Region**

**Application storefront**

Purchasing, subscriptions, entitlements, catalog → Storefront services → Client Edge Living Room, Mobile, Web → Network edge

**Playback**

DRM, Bookmarking, Advertising → Playback services → Playback URLs → Playback Edge

**Live signal delivery**

Production → Signal Acquisition → Encoding & Packaging → Content Management & Delivery

Live Source

Customers

Content Delivery Networks

# 10,000 ft architecture



AWS Region

## Application storefront

Purchasing, subscriptions, entitlements, catalog → Storefront services → Client Edge Living Room, Mobile, Web → Network edge

## Playback

DRM, Bookmarking, Advertising → Playback services → Playback URLs → Playback Edge

## Live signal delivery

Live Source → Production → Signal Acquisition → Encoding & Packaging → Content Management & Delivery

Customers

Content Delivery Networks

# 10,000 ft architecture



**AWS Region**

**Application storefront**
Purchasing, subscriptions, entitlements, catalog — Storefront services — Client Edge Living Room, Mobile, Web — Network edge

**Playback**
DRM, Bookmarking, Advertising — Playback services — Playback URLs — Playback Edge

**Live signal delivery**
Production — Signal Acquisition — Encoding & Packaging — Content Management & Delivery

Live Source

Customers

Content Delivery Networks

# 10,000 ft architecture



AWS Region

## Application storefront

Purchasing, subscriptions, entitlements, catalog — Storefront services — Client Edge Living Room, Mobile, Web — Network edge

## Playback

DRM, Bookmarking, Advertising — Playback services — Playback URLs — Playback Edge

## Live signal delivery

Live Source — Production — Signal Acquisition — Encoding & Packaging — Content Management & Delivery

Content Delivery Networks

Customers

# Signal ingestion: Stadium source

# Signal ingestion: Mobile production studio

# Signal ingestion: Eyes on glass

# Signal ingestion: Architecture deep dive



Venue

Production & media distribution

Ad playout → Encoders

Ad playout → Encoders

AWS Direct Connect

Edge reception and routing

AWS Elemental Media Stack

| MediaConnect | MediaLive | MediaPackage |

Origin1: us-east-1

| MediaConnect | MediaLive | MediaPackage |

Origin2: us-west-2

| MediaConnect | MediaLive | MediaPackage |

Origin3: us-east-1

| MediaConnect | MediaLive | MediaPackage |

Origin4: us-west-2

CDN 1

CDN 2

CDN 3

CDN 4

CDN 5

CDN 6

# Signal ingestion: Architecture deep dive



Venue

Production & media distribution

Ad playout → Encoders

Ad playout → Encoders

AWS Direct Connect

Edge reception and routing

AWS Elemental Media Stack

| MediaConnect | MediaLive | MediaPackage |
Origin1: us-east-1

| MediaConnect | MediaLive | MediaPackage |
Origin2: us-west-2

| MediaConnect | MediaLive | MediaPackage |
Origin3: us-east-1

| MediaConnect | MediaLive | MediaPackage |
Origin4: us-west-2

CDN 1
CDN 2
CDN 3
CDN 4
CDN 5
CDN 6

# 10,000 ft architecture

# Multi-Region playback

**Enhanced availability and resilience
for NFL Thursday Night Football**

**Regional failover**

**Edge routing (PortKey)**

**Active/active**

**Globalization**

# Multi-Region architecture



Routing rules

Routing rules

Network edge services

Routing layer

Web/client edge

Region 1

Client

Network edge services

Routing rules

Routing layer

Web/client edge

Region 2

Network edge services

Routing rules

Routing layer

Web/client edge

Region 3

# Multi-Region architecture



Routing rules

Routing rules

Routing layer

Network edge services

Web/client edge

Region 1

Client

Routing rules

Routing layer

Network edge services

Web/client edge

Region 2

Routing rules

Routing layer

Network edge services

Web/client edge

Region 3

# Multi-Region architecture



Region 1

Region 2

Region 3

Client

Network edge services

Routing rules

Routing layer

Web/client edge

# Multi-Region architecture



Routing rules

Routing rules

Routing layer

Web/client edge

Region 1

Network edge services

Client

Network edge services

Routing rules

Routing layer

Web/client edge

Region 2

Network edge services

Routing rules

Routing layer

Web/client edge

Region 3

# Multi-Region architecture



Routing rules

Routing rules

Network edge services

Routing layer

Web/client edge

Region 1

Client

Network edge services

Routing rules

Routing layer

Web/client edge

Region 2

Network edge services

Routing rules

Routing layer

Web/client edge

Region 3

# Key benefits of multi-Region playback

**Enhanced availability and resilience
for NFL Thursday Night Football**

**Lower latency**

**Enhanced resilience**

**Seamless failover**

# 10,000 ft architecture

# Multi-Region storefront and detail pages

- App start, home page, and detail page
- Globalization: Feature parity between Regions
- Data replication and consistency
- Shadow testing
- PhasedlLaunch

# Data replication and consistency

- Leveraged Amazon DynamoDB global tables for low-latency access to critical metadata
- Bookmarking multi-master cross-region replication strategy

# 10,000 ft architecture



**AWS Region**

**Application storefront**

Purchasing, subscriptions, entitlements, catalog → Storefront services → Client Edge Living Room, Mobile, Web → Network edge

**Playback**

DRM, Bookmarking, Advertising → Playback services → Playback URLs → Playback Edge

**Live signal delivery**

Production → Signal Acquisition → Encoding & Packaging → Content Management & Delivery

Live Source

Customers

Content Delivery Networks

# Real-world success stories

- In 2024 we served 50% of customers from another Region
- Millions of concurrent users
- Smooth user experiences despite high demand surges

# Lessons learned and future enhancements

**Enhanced availability and resilience
for NFL Thursday Night Football**

**Automation**

**Testing**

**Always-on**

**Iteration**

An elastic scaling story

# TRIVIA

**In 2024, what was the peak audience for the Cowboys-Giants Thursday Night Football game on Prime Video?**

A. 18 million

B. 10 million

C. 15 million

D. 20 million

# TRIVIA

**In 2024, what was the peak audience for the Cowboys-Giants Thursday Night Football game on Prime Video?**

**A.** 18 million

**B.** 10 million

**C.** 15 million

**D.** 20 million

# Scaling for Live Sports at Prime Video

- Impact of live sports
- Prime Video's role in streaming Thursday Night Football (TNF)

**Peak Concurrency**

Series1  Series2

# Scaling for Live Sports at Prime Video

- Impact of live sports
- Prime Video's role in streaming Thursday Night Football (TNF)



Peak Concurrency

Series1 ——— Series2

# Scaling for Live Sports at Prime Video

- Impact of live sports
- Prime Video's role in streaming Thursday Night Football (TNF)

# Key challenges to scaling for peak NFL events

- Overall scale
- Spiky traffic patterns
- Coordination Effort & Lead time
- Lack of Automation

**Peak Concurrency**

| 1-Jan | 8-Jan | 15-Jan | 22-Jan | 29-Jan |

Series1    Series2

# Key challenges to scaling for peak NFL events

- Overall scale
- Spiky traffic patterns
- Coordination Effort & Lead time
- Lack of Automation

**TNF traffic ramp**

**Halftime**

1  3  5  7  9  11131517192123252729313335373941434547495153555759

# Key challenges to scaling for peak NFL events

- Overall scale
- Spiky traffic patterns
- Coordination Effort & Lead time
- Lack of Automation

# Key challenges to scaling for peak NFL events

- Overall scale
- Spiky traffic patterns
- Coordination Effort & Lead time
- Lack of Automation

**Peak Concurrency**

1-Jan    8-Jan    15-Jan    22-Jan    29-Jan

Series1    Series2

# Opportunity with auto scaling (elasticity)

- Address the area under the curve

- Reduce number of days at peak scale

- Reduce infra $$

- Reduce engineering effort

## Peak Concurrency

1-Jan         8-Jan         15-Jan         22-Jan         29-Jan

Series1      Series2

Scaling Waste

# Signaling solution to enable auto scaling

- Proactive fully-automated signaling solution
- Building a solution that is AWS product-agnostic, and can be leveraged across hundreds of teams

# Prime Video automated signaling solution



Forecasting AWS Accounts

Prime Video Forecasting Solution

Notifications Amazon SNS

Ahead-of-time Solution AWS Accounts

Encryption AWS KMS

Data Bucket Amazon S3

Amazon EventBridge

Public Certificate AWS Certificate Manager

Domains Amazon Route 53

Notification Collector Amazon ECS

Retry Mechanism Step Function

Service Amazon ECS

Application Load Balancer

Notification Queue Amazon SQS

Queue Depth Metric Amazon CloudWatch

Amazon ECS Desired Count App Auto Scaling

Task Scheduler Amazon ECS

Task Data Amazon DynamoDB

Prime Video Service AWS Account

EventBridge

Agent AWS Lambda

App Auto Scaling

Customer Service Amazon ECS/ Amazon EC2

Metrics CloudWatch

aws

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Proactive demand forecasting

- Utilizing machine learning and business inputs to predict NFL demand

- Importance of accuracy in capacity planning

**Peak Concurrency**

——Series1

# Prime Video automated signaling solution

# Prime Video automated signaling solution
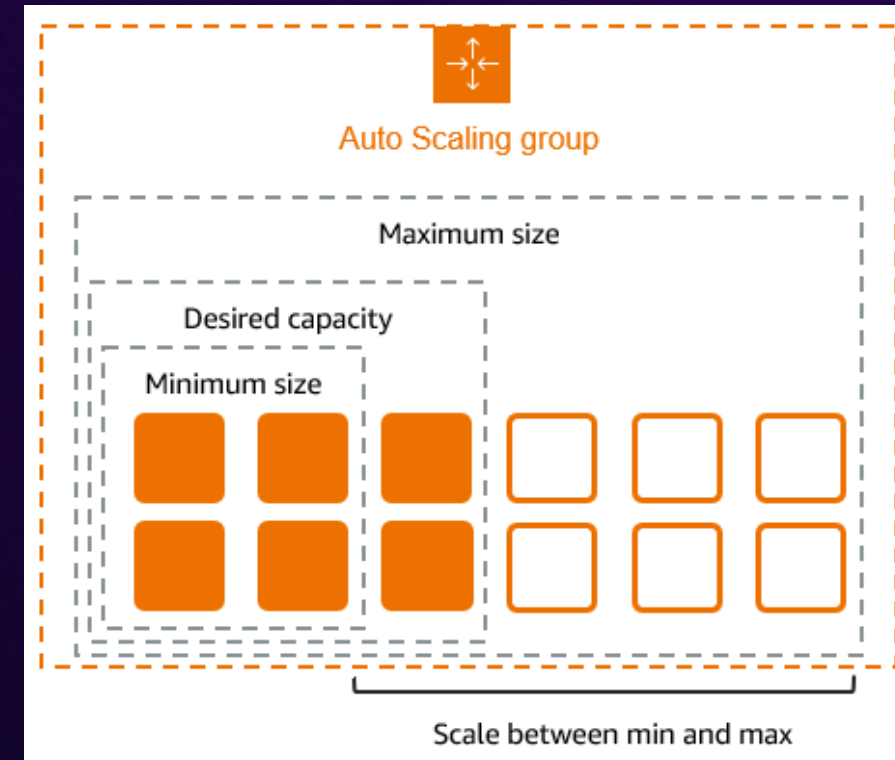
# Scheduled (proactive) Auto Scaling

# Scheduled (proactive) Auto Scaling

# Scheduled (proactive) Auto Scaling

## CloudWatch Metric Dashboard

# Dynamic (reactive) auto scaling

- Real-time scaling to handle unanticipated spikes in traffic:
  - Popular Cowboys matchup exceeding forecasts
  - Taylor Swift surprise appearance!
- Benefits of immediate responsiveness
  - Monitors key metrics and scales out resources within minutes

# Success Metrics and Impact - Pilot

Success metrics for evaluating auto scaling performance for pilot services

- Reduction in days at peak scale of ~70%
- Improvements in CPU utilization by ~20-30%
- Cost reduction of ~30%

# Lessons learned and best practices

- Understand your workloads (scheduled vs. dynamic)
- Complement with dynamic auto scaling
- Start simple, gradually increase
- Monitoring, Alarms, and optimize KPIs

# Future of Prime Video



2016                    2018                    2020                    2022

# Recap

# Recap

**1**

**Delivery – Premium video at scale to millions of customers in across thousands of different devices**

**2**

**Customer experience – Millions of requests and terabits of data over highly variable networks**

**3**

**Availability and resiliency – Achieve 100% up-time with multi-region and autoscaling strategies**

# Number of AWS services?

# 120+

aws

# AWS Well-Architected Framework



Operational excellence

Security

Reliability

Cost optimization

Sustainability

Performance efficiency

Reliability

# Reliability (Goal)

**The ability of a system to perform its required function correctly and consistently**
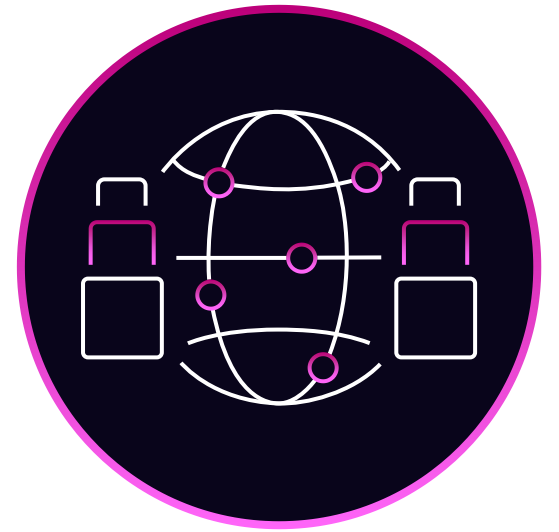
# Resilience (Approach)

**The ability of a system to recover quickly and restore to a fully functional state after a disruption or failure**

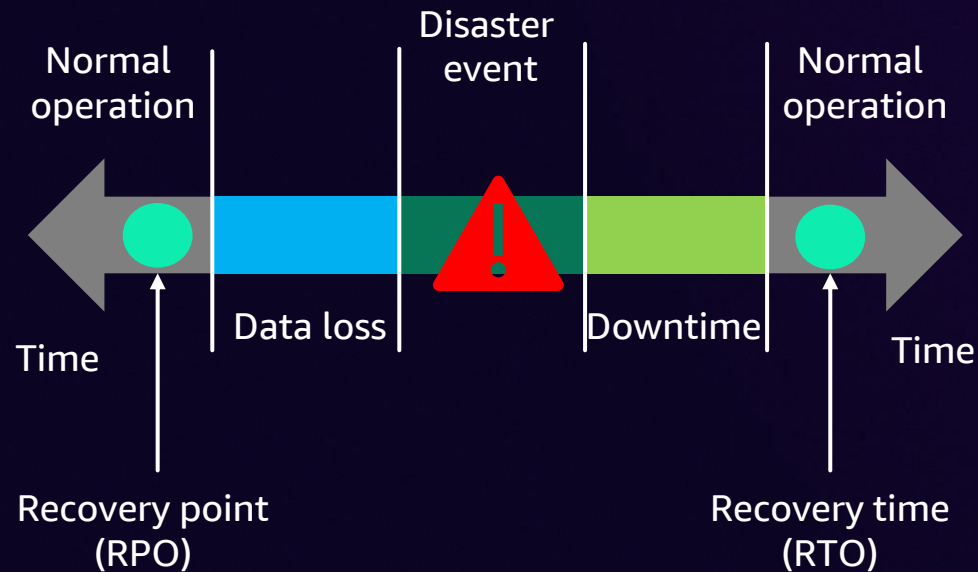**https://bit.ly/reliability-pillar**

# Reliability design principles

- Automatically recover from failure

---

- Test recovery procedures

---

- Scale horizontally to increase aggregate workload availability

---

- Stop guessing capacity

---

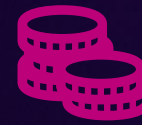- Manage change through automation



aws

# Multi-Region best practices

- Understand consistency and replication options

- Utilize asynchronous writes across Regions

- Multi-Region writes should be idempotent

- Observability from multiple viewpoints

# Fundamental #1: Understand the requirements



Normal operation — Disaster event — Normal operation

Data loss | Downtime

Time — Time

Recovery point (RPO)

Recovery time (RTO)

Why multi-Region?
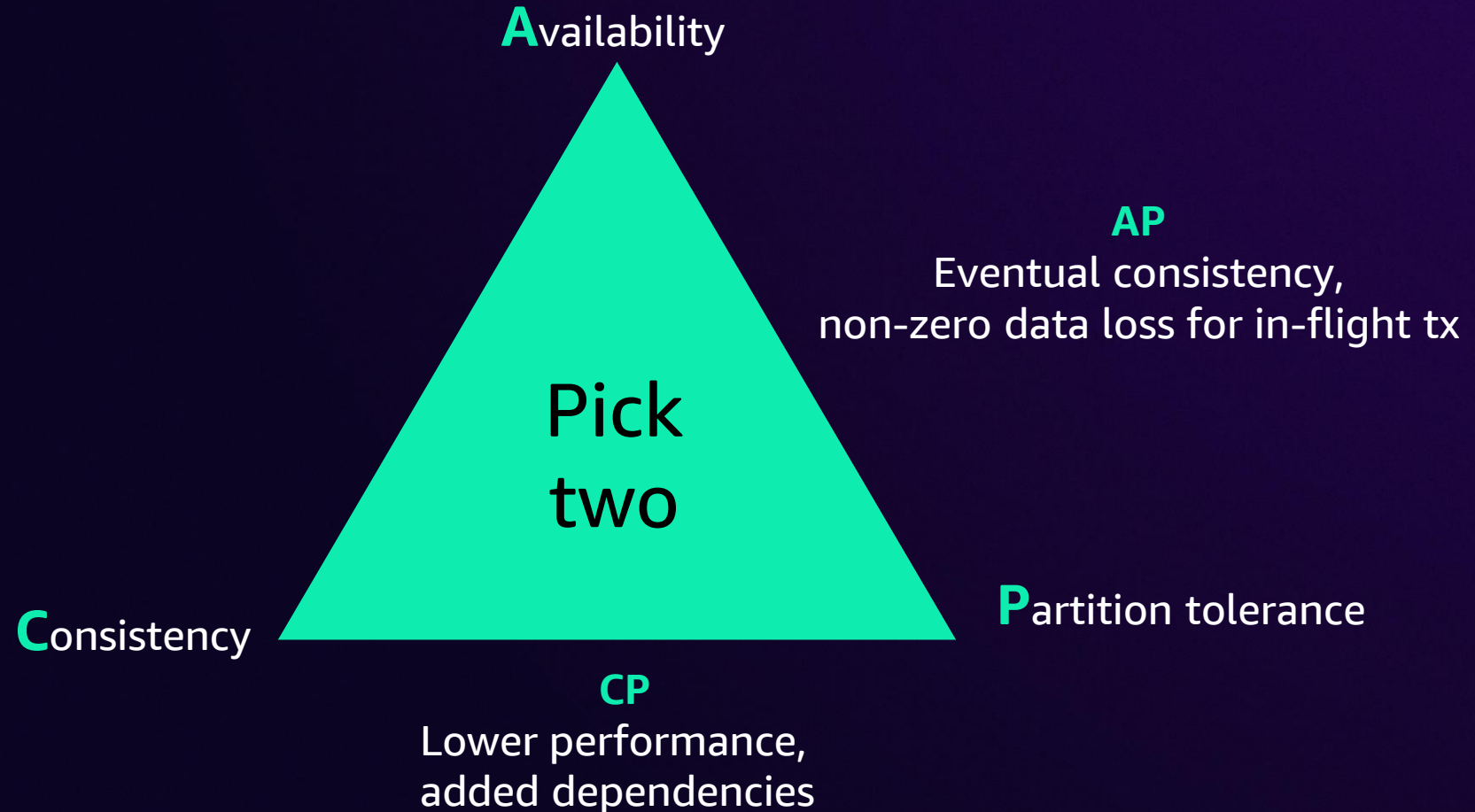
Cost to run your application multi-Region vs. business impact

What's the RTO/RPO requirements?

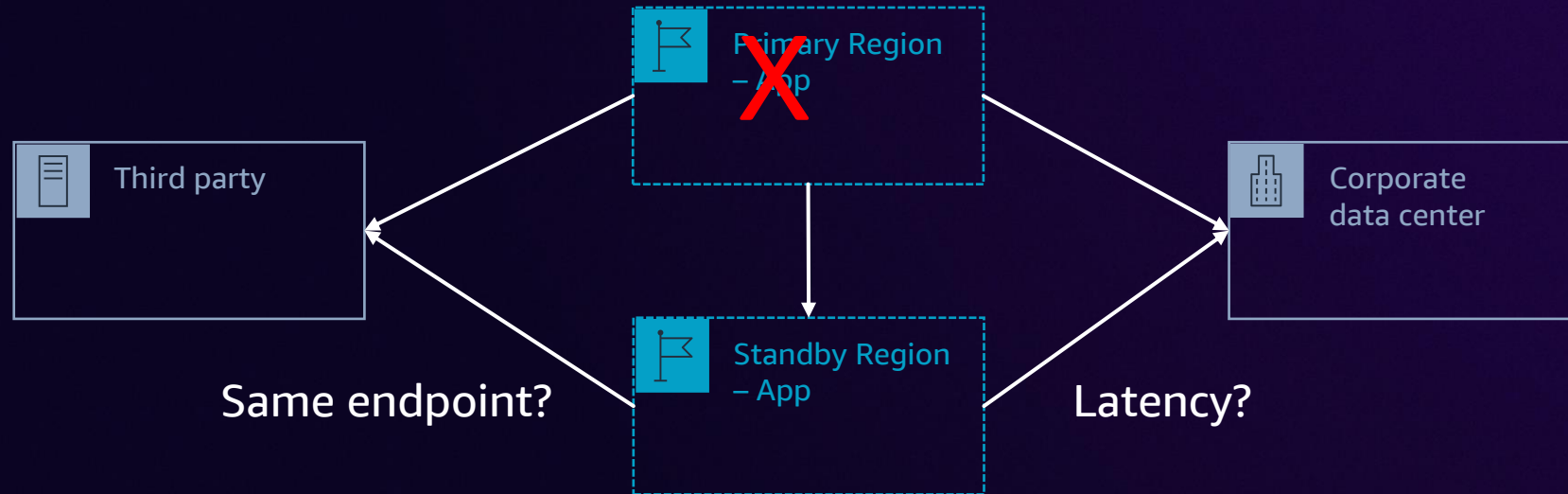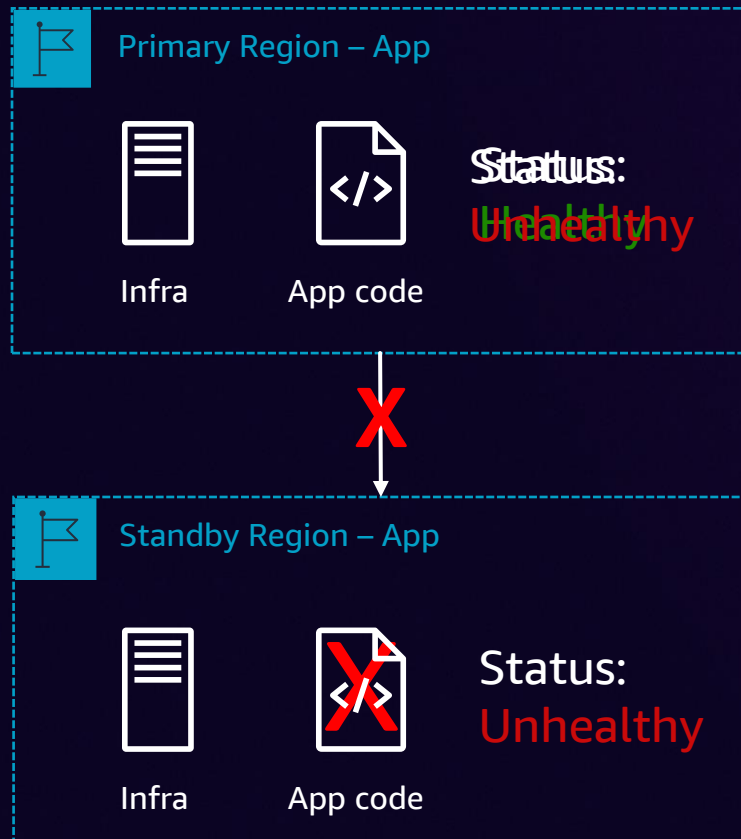Development cycles to improve application resilience?

# Fundamental #2: Understand data consistency

**A**vailability

**Pick two**

**AP**
Eventual consistency,
non-zero data loss for in-flight tx

**C**onsistency

**P**artition tolerance

**CP**
Lower performance,
added dependencies

# Fundamental #3: Understanding dependencies



Third party

Primary Region – App

Standby Region – App

Corporate data center

Same endpoint?
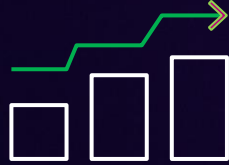
Latency?

# Key takeaways

# Key takeaways

**Operational resilience**

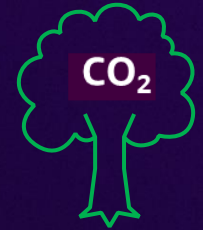Support growth for millions of users on TNF days

**Dynamic scaling**

Peak-to-mean ratio is **highly variable** (baseline days to NFL days)

**Optimized cost**

**~20%** reduction in cost and operations

**Carbon savings**

**43%** reduction in carbon footprint by reducing number of EC2 instances

Resiliency →→→ Optimization

"

**Everything fails all the time. We needed to build systems that embrace failure as a natural occurrence.**

**Dr. Werner Vogels**

CTO, Amazon

aws

# Thank you!

Please complete the session survey in the mobile app

**Tulip Gupta**
✉ tulipgpt@amazon.com
in linkedin.com/in/tulip-gupta

**Elliott Nash**
✉ ellinash@amazon.com
in linkedin.com/in/elliottnash

**Ralph Chaker**
✉ chakerr@amazon.com
in linkedin.com/in/ralphchaker