

ARC330

AWS re:INVENT

How the BBC Built a Massive Media Pipeline Using Microservices

Stephen Godwin, Lead Architect, British Broadcasting Corporation (BBC)
Robert Chow, Principal Product Manager, AWS (SQS and SNS)

November 28, 2017

AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



The BBC iPlayer is the biggest audio and video-on-demand service in the UK. Over one-third of the country submits 10 million video playback requests every day, and the service publishes over 10,000 hours of media every week. Moving iPlayer to the cloud has enabled the BBC to shorten the time-to-market of content from 10 hours to 15 minutes. In this session, the BBC's lead architect describes the approach behind creating iPlayer architecture, which uses **Amazon SQS** and **Amazon SNS** in several ways to improve elasticity, reliability, and maintainability. You see how BBC uses AWS messaging to choreograph the 200 microservices in the iPlayer pipeline, maintain data consistency as media traverses the pipeline, and refresh caches to ensure timely delivery of media to users. This is a rare opportunity to see the internal workings and best practices of one of the largest on-demand content delivery systems operating today.

Where we're going

Messaging overview

BBC iPlayer

- Background
- Challenges
- Transformation to microservices
- The role of messaging
- Results
- Pitfalls and best practices

Wrap-up

AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



System building blocks for modern applications



Compute



Storage



Database

AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



System building blocks for modern applications



Compute



Storage



Database



Messaging

AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



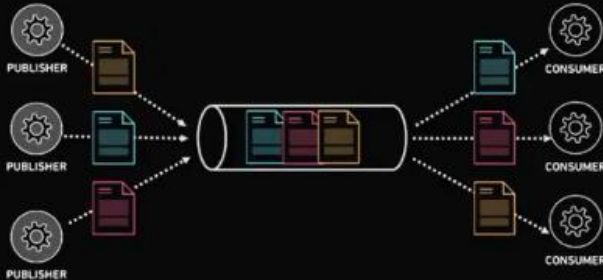
We also think about messaging using 2 specific services **SQS** and **SNS** that are absolute critical for data management in a scalable system

Amazon SQS & SNS: Decouple applications



Simple Queue Service (SQS)

- Asynchronous
- Each message processed by one consumer



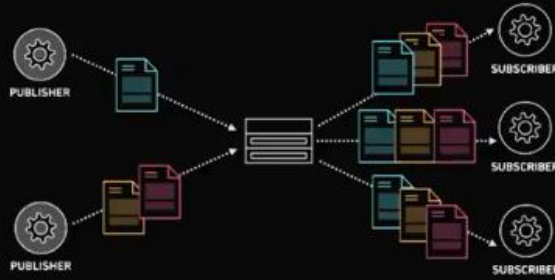
AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



Simple Notification Service (SNS)

- Pub/sub
- Broadcast or filtered
- High fan-out of messages



aws

Stephen Godwin **BBC**



@SteveGodwin

- Lead technical architect at the BBC
- Designed the systems that provides audio and video to BBC iPlayer and iPlayer Radio
- Led the migration of the systems that power iPlayer to a cloud-based microservice architecture
- Worked for nearly 10 years as one of the developers of IBM WebSphere

AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

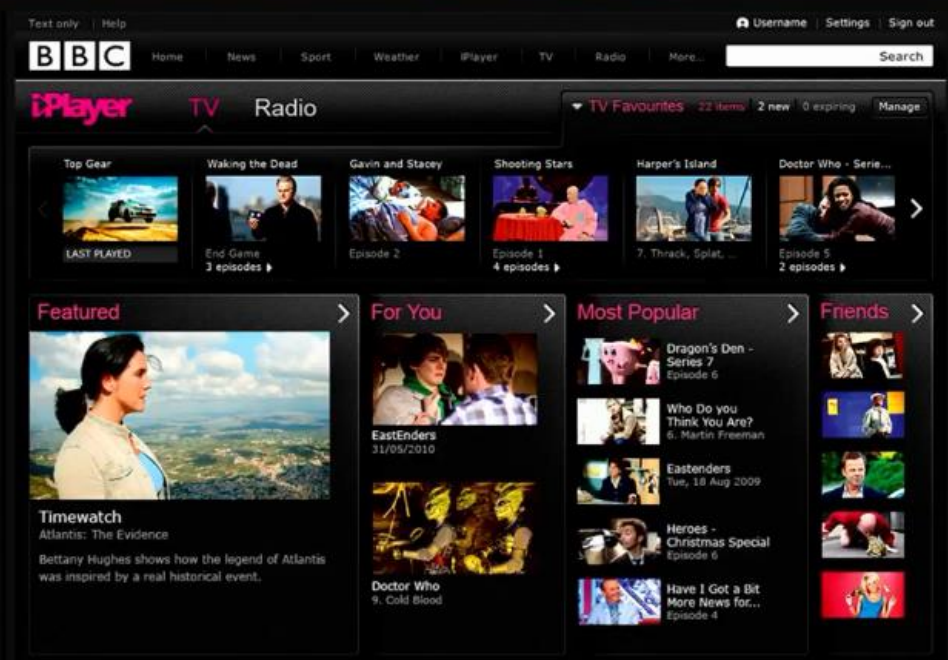
aws

BBC



AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



BBC iPlayer

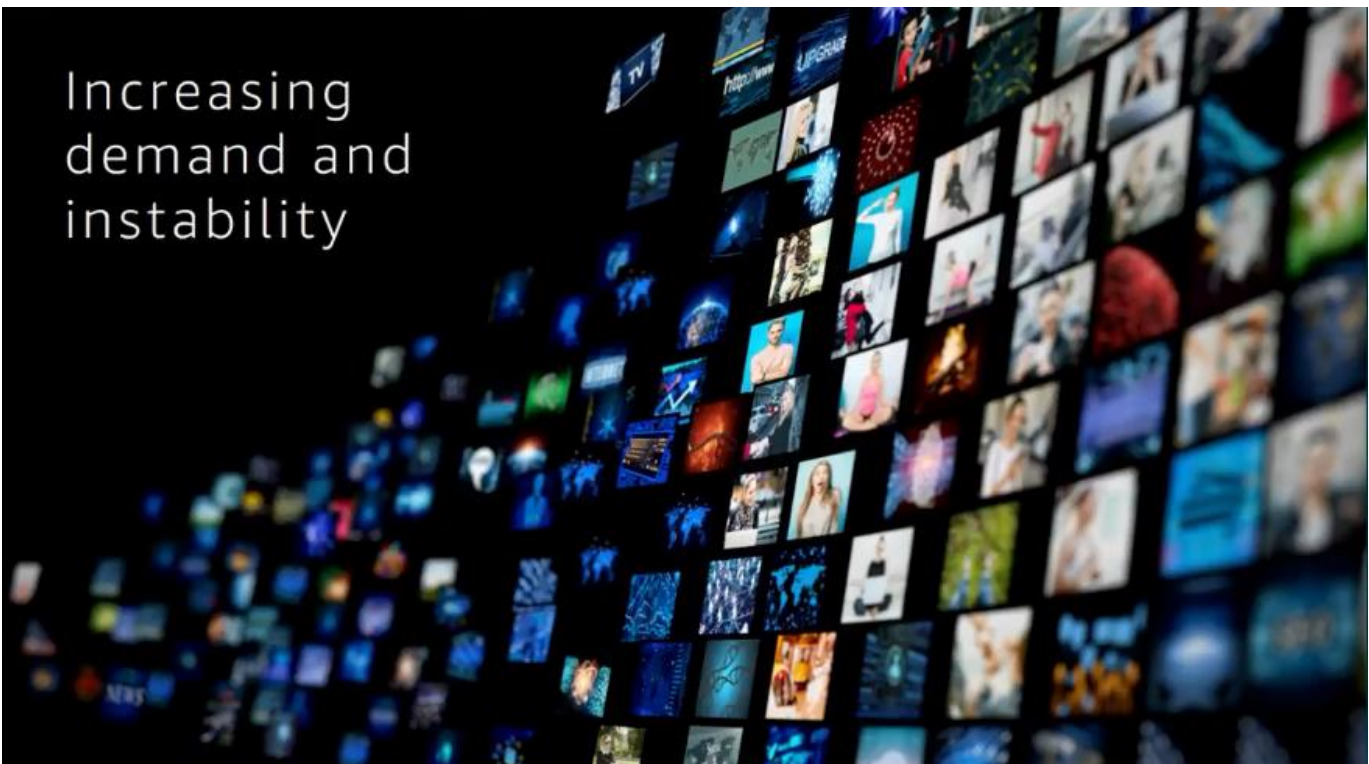
WATCHED BY

30 PERCENT
OF THE UK
POPULATION

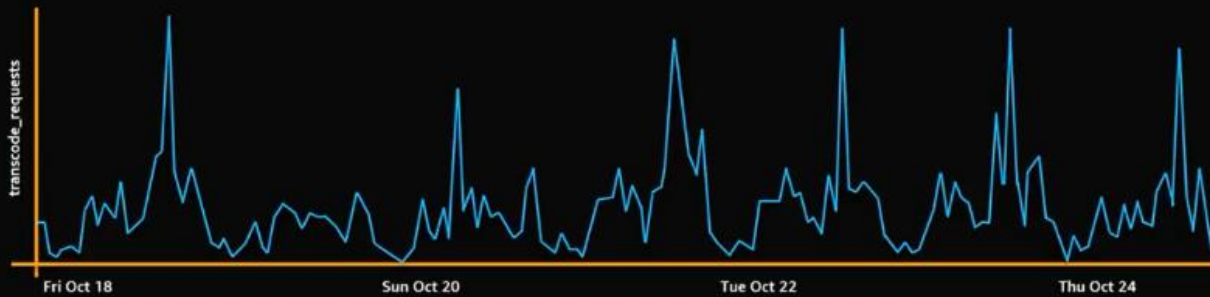
10 THOUSAND HOURS
OF MEDIA PUBLISHED
EVERY WEEK

EVERY
A NEW PROGRAM
IS AVAILABLE
ON iPLAYER
15
MINUTES

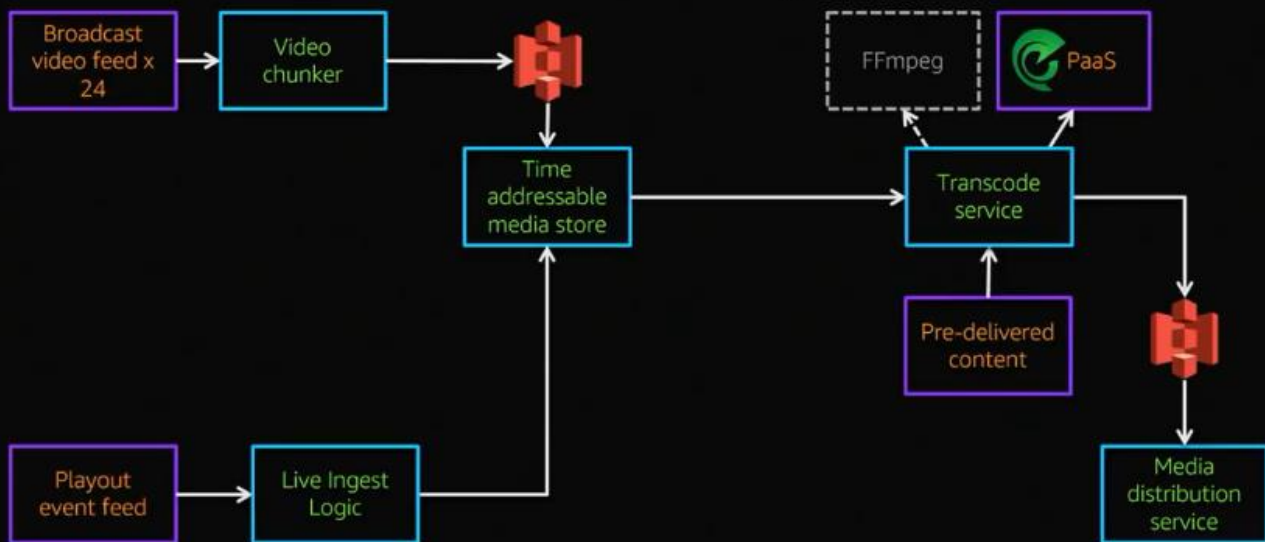
Increasing
demand and
instability



A perfect pattern for an elastic architecture



Inside iPlayer: Video factory

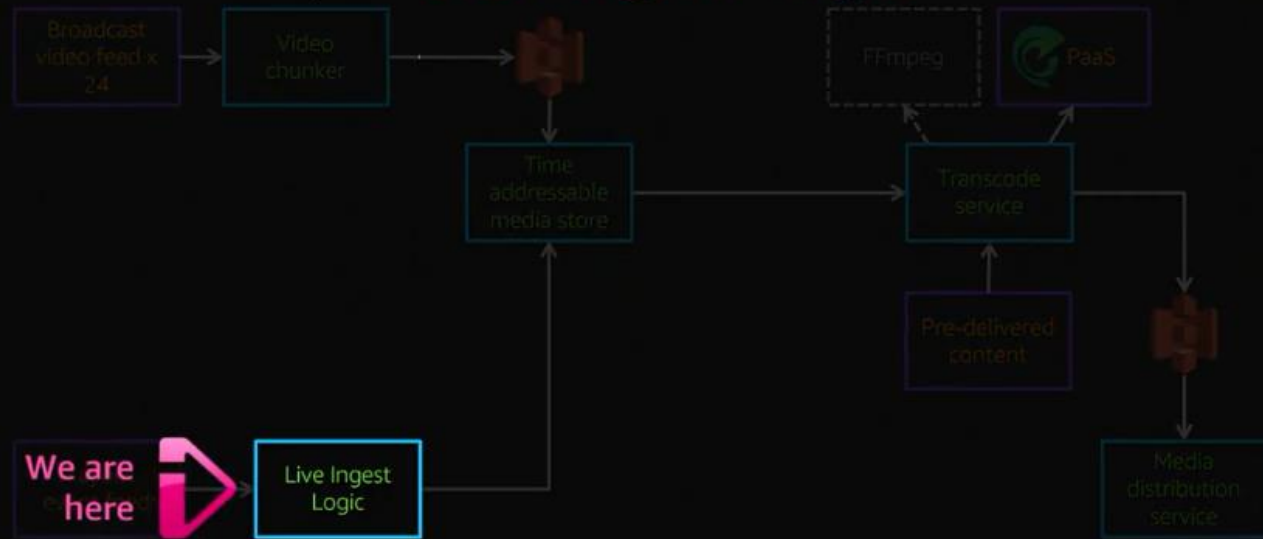


AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



Inside iPlayer: Live Ingest Logic



AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

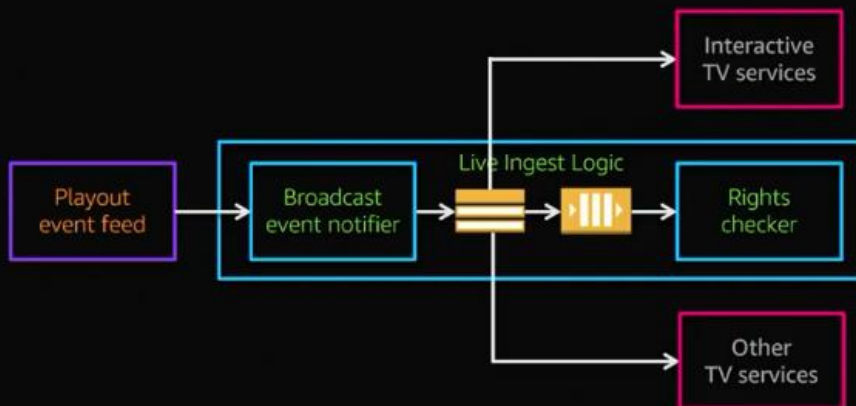


Inside the Live Ingest Logic service

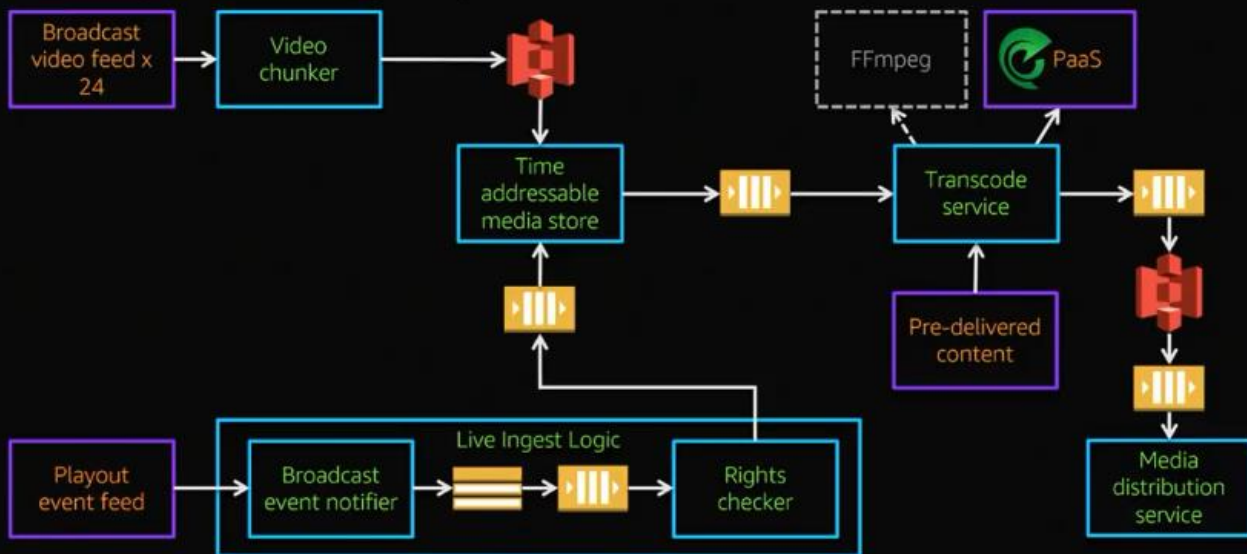


This is a pattern where an SQS queue subscribes to an SNS topic like startEvent or stopEvent

Inside the Live Ingest Logic service



Using messaging to wire the system

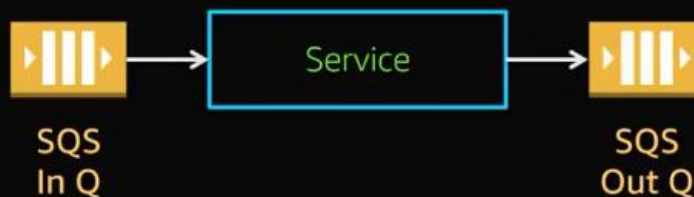


AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



Message-oriented microservices pattern

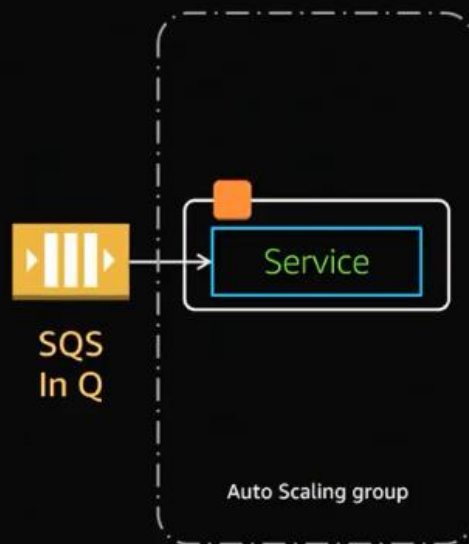


Message-oriented microservices pattern

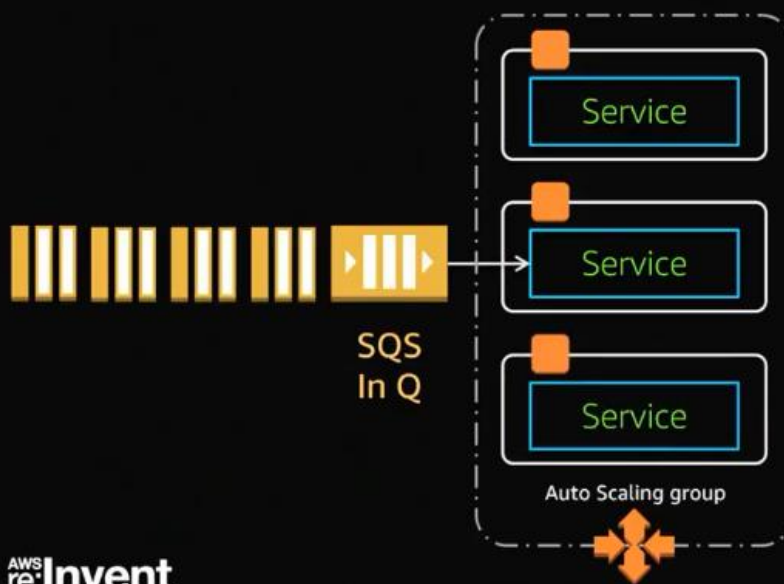


Running on each Service is a Java application running on an EC2 instance that doesn't know much about the queues when doing its job.

Scaling to handle peaks



Scaling to handle peaks



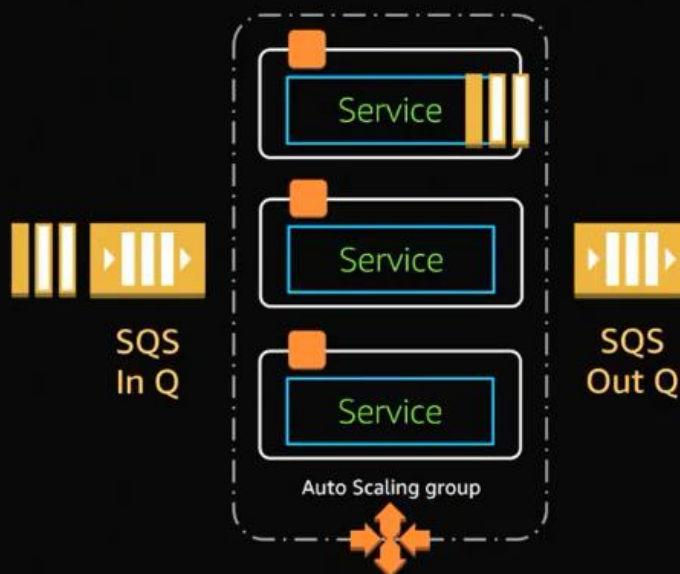
aws re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

aws

We normally run at least 3 instances in a different auto-scaling group as above so as to scale up or down when needed.

Resiliency and durability

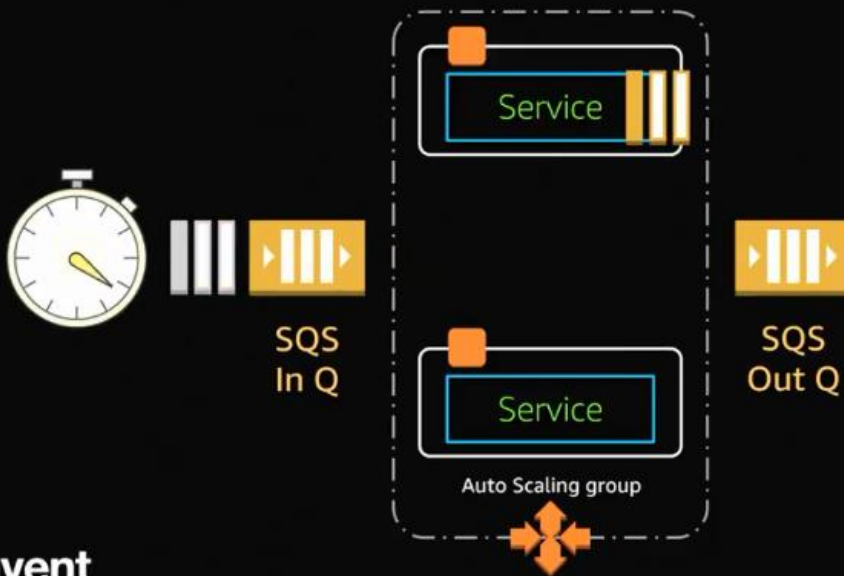


aws re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

aws

Resiliency and durability



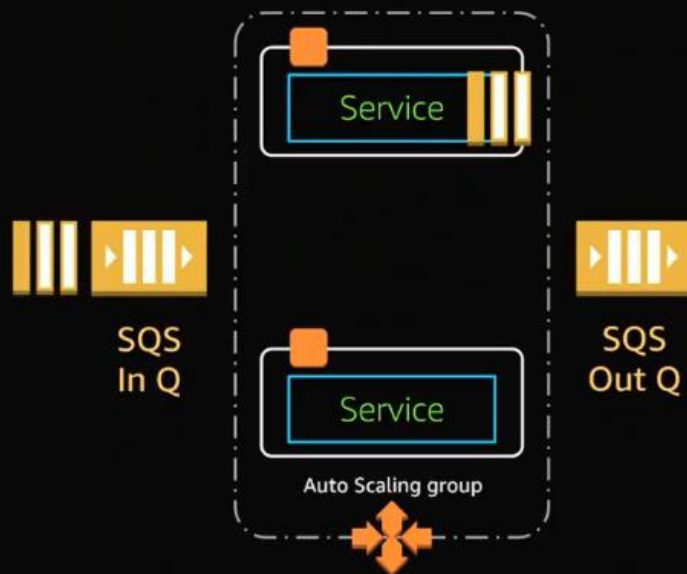
AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

aws

If a message processing fails due to an EC2 instance getting terminated or faulty, then the message will reappear back in the SQS In Q after a certain period of time (Visibility Timeout) if it hasn't generated a message onto the SQS Out Q queue

Resiliency and durability

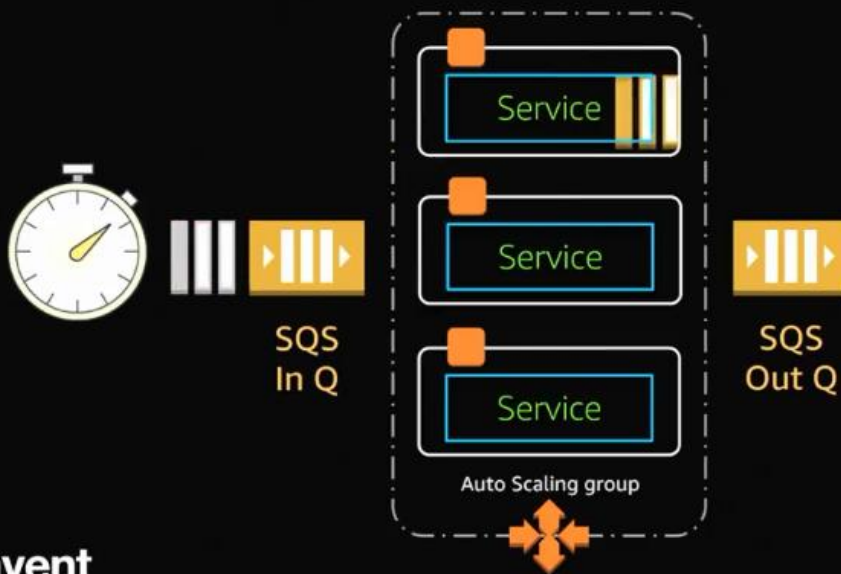


AWS re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

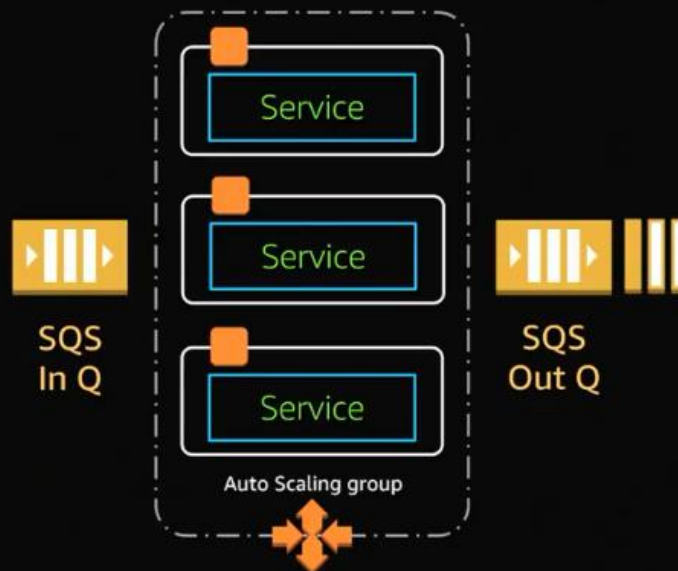
aws

Resiliency and durability



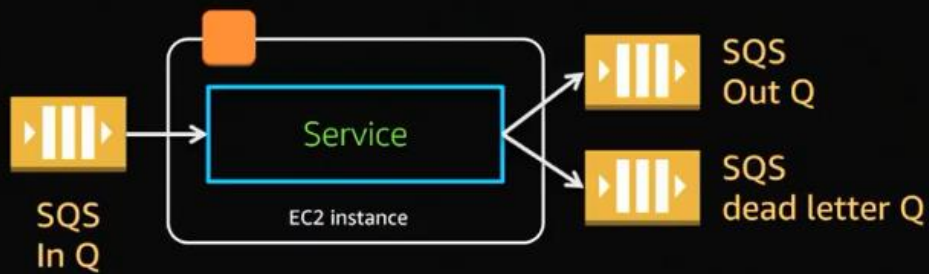
This is what happens in normal processing, an EC2 instance takes a message off the queue for processing and passing the result as a message in the Out Q queue as below

Resiliency and durability

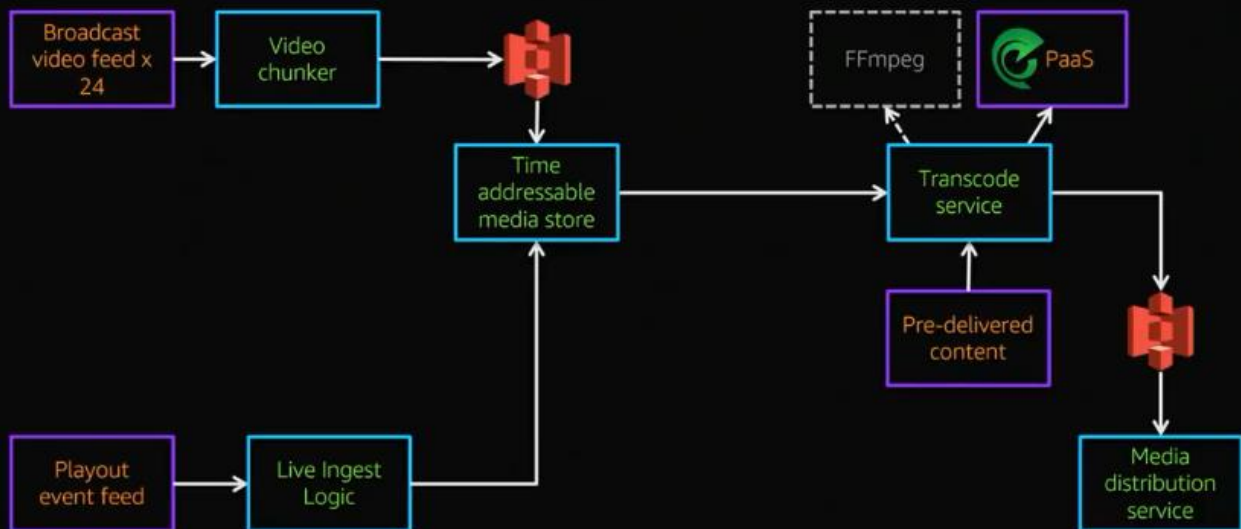


Only after the result message has been placed in the SQS Out Q do we then delete the equivalent message from the SQS In Q queue

Error handling



Inside iPlayer: Time addressable media store

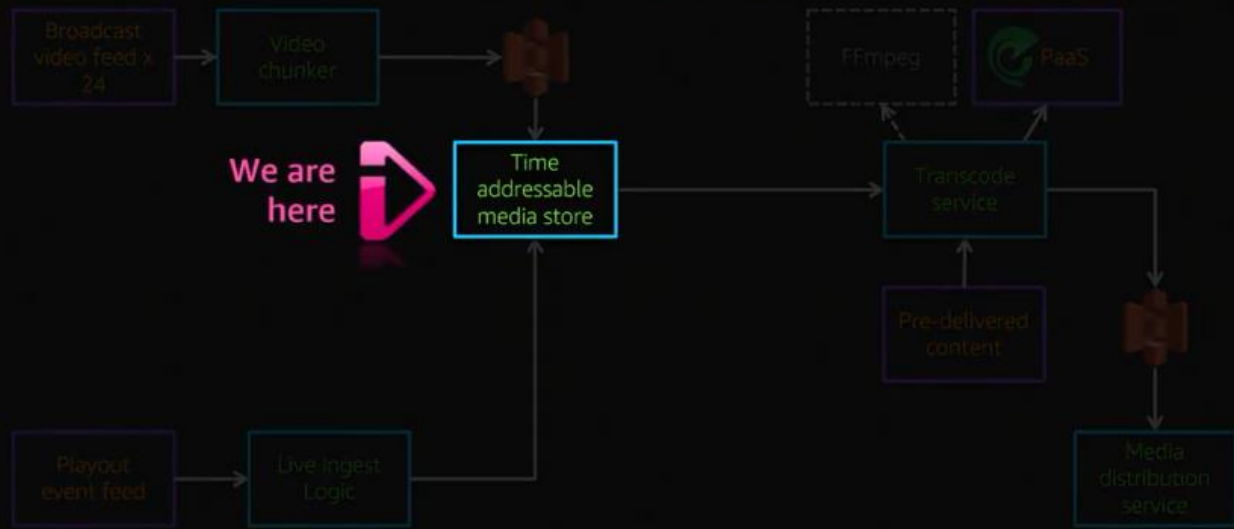


AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

aws

Inside iPlayer: Time addressable media store

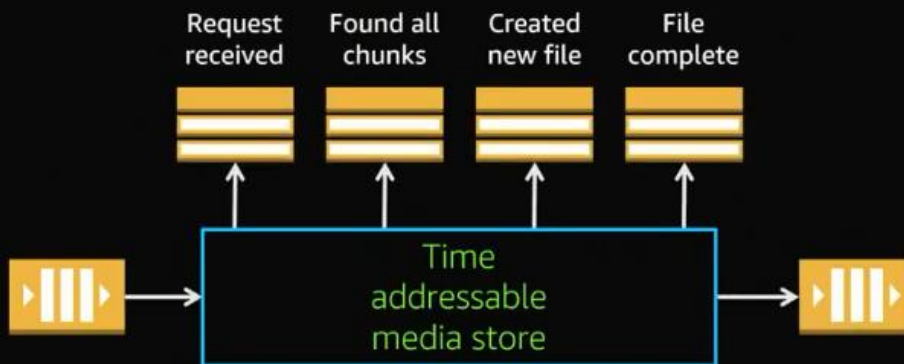


AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

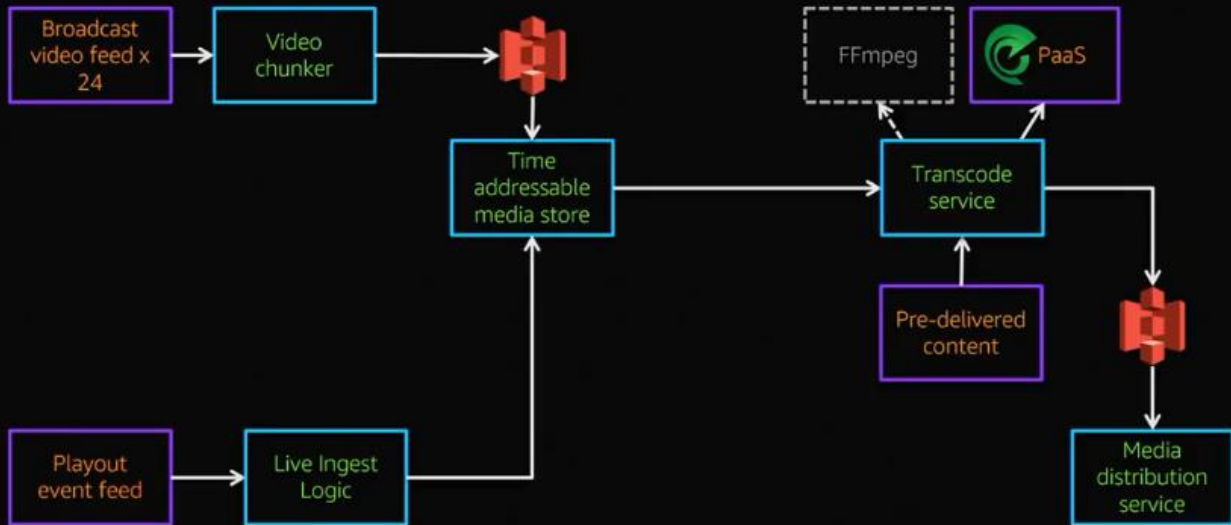


Event-driven trace and analysis



Each message entering the system gets given **a unique correlation identifier** that gets copied from service to service as it moves through the processing pipeline, this allows us to know what event logs are from what message. We can use a particular correlation ID to see all events or messages have been triggered from that initial message

Inside iPlayer: Media distribution service



AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



Inside iPlayer: Media distribution service



We are all the way over here

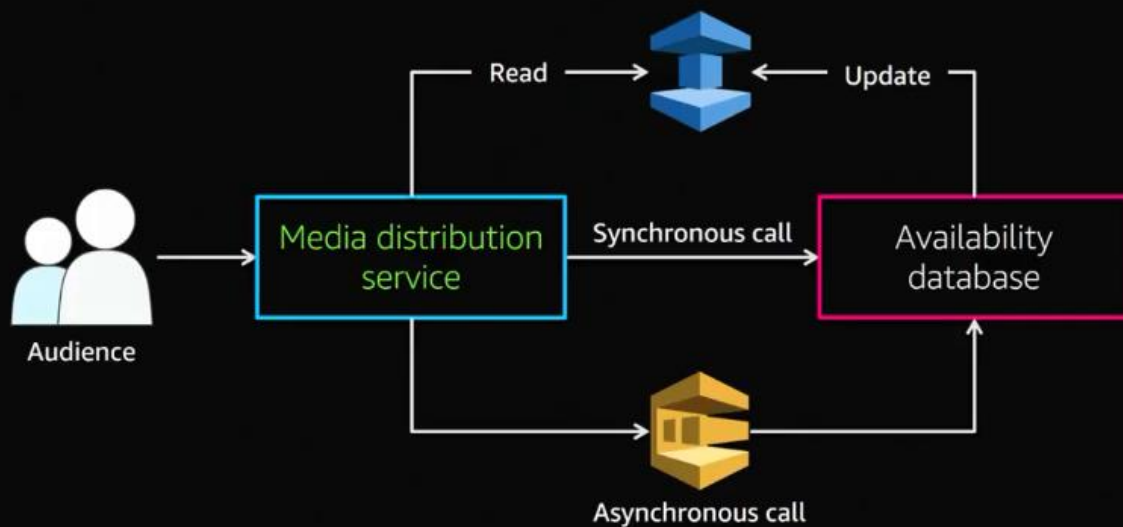


AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



Refreshing caches for data consistency



AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



Every time a user tries to view a video content, this service gets called about what **CDN** should be used to get the nearest copy of the video or file by querying the **Availability database**. We also cache the result in **AWS ElasticCache** so that we don't need to query again for this particular user and content. If a user makes a video request and we query the cache and get a hit, if the cache period for that content is going to be stale in the next 1 minute, we will put in an async request into the queue to refresh that content when possible and reset the TTL in the cache for continuous content freshness.

Where we ended up

Happier audience

Faster time to deployment

Easy extensibility and reuse

Massive scalability

Improved resilience

AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



LESSONS LEARNED

Designing for idempotency
Elastic scaling is great—linear scaling is better
Reliable pipelines of microservices
Think about how you will debug your system

what's next



LEARN MORE

aws.amazon.com/sqs/getting-started

aws.amazon.com/sns/getting-started

AWS
re:Invent

© 2017, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

