# Polyglot Persistence Powering Microservices

Roopa Tangirala

Engineering Manager

Netflix

# Agenda

- 5 Use Cases
- Challenges
- Current Approach
- Takeaway

# About Netflix

Netflix has been leading the way for digital content since 1997

## NETFLIX ORIGINALS

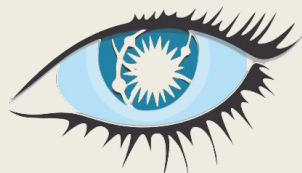NETFLIX ORIGINAL
**STRANGER THINGS**
NEW EPISODES

NETFLIX ORIGINAL
**ZUMBO'S JUST DESSERTS**

NETFLIX ORIGINAL
**ALIAS GRACE**

NETFLIX ORIGINAL
JACK WHITEHALL:
**TRAVELS WITH MY FATHER**

NETFLIX ORIGINAL
**THE SINNER**

NETFLIX ORIGINAL
**6 DAYS**

## Trending Now

**GET HARD**

NETFLIX
**NARCOS**

NETFLIX
**DYNASTY**
NEW EPISODE    WEEKLY

**POWER**

**HOMELAND**

**SUITS**

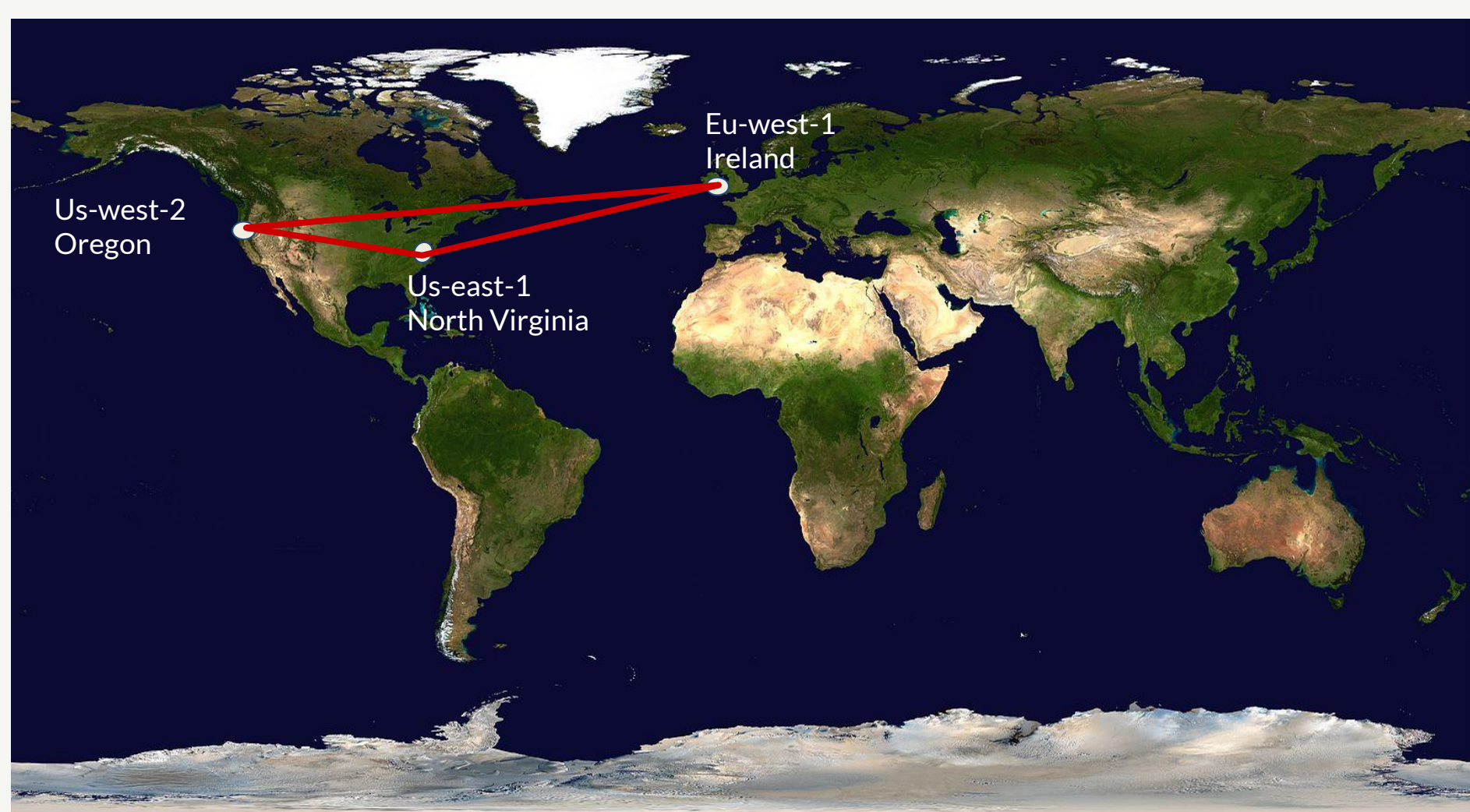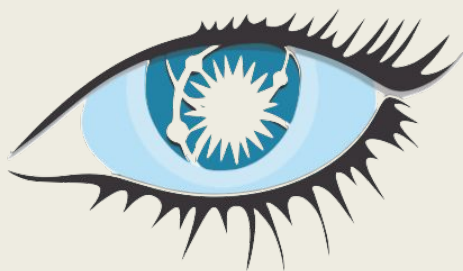| | |
|---|---|
|  elastic | Search, Analyze and visualize in near real time |
|  EVCache | Distributed in-memory caching solution based on memcached |
|  cassandra | Distributed NOSQL database to handle large datasets providing high availability. |
|  DYNOMITE | Distributed dynamo layer for different storage engines and protocols supporting Redis, memcached, RocksDB |
|  | TitanDB is scalable graph database optimized for storing and querying graph datasets. |

# Requirements - CDN URL

- High availability
- Very low latency reads/writes (less than 1ms)
- High Throughput per node

# ★ Playback Error



Whoops, something went wrong...

**Unexpected Error**

There was an unexpected error. Please reload the page and try again.

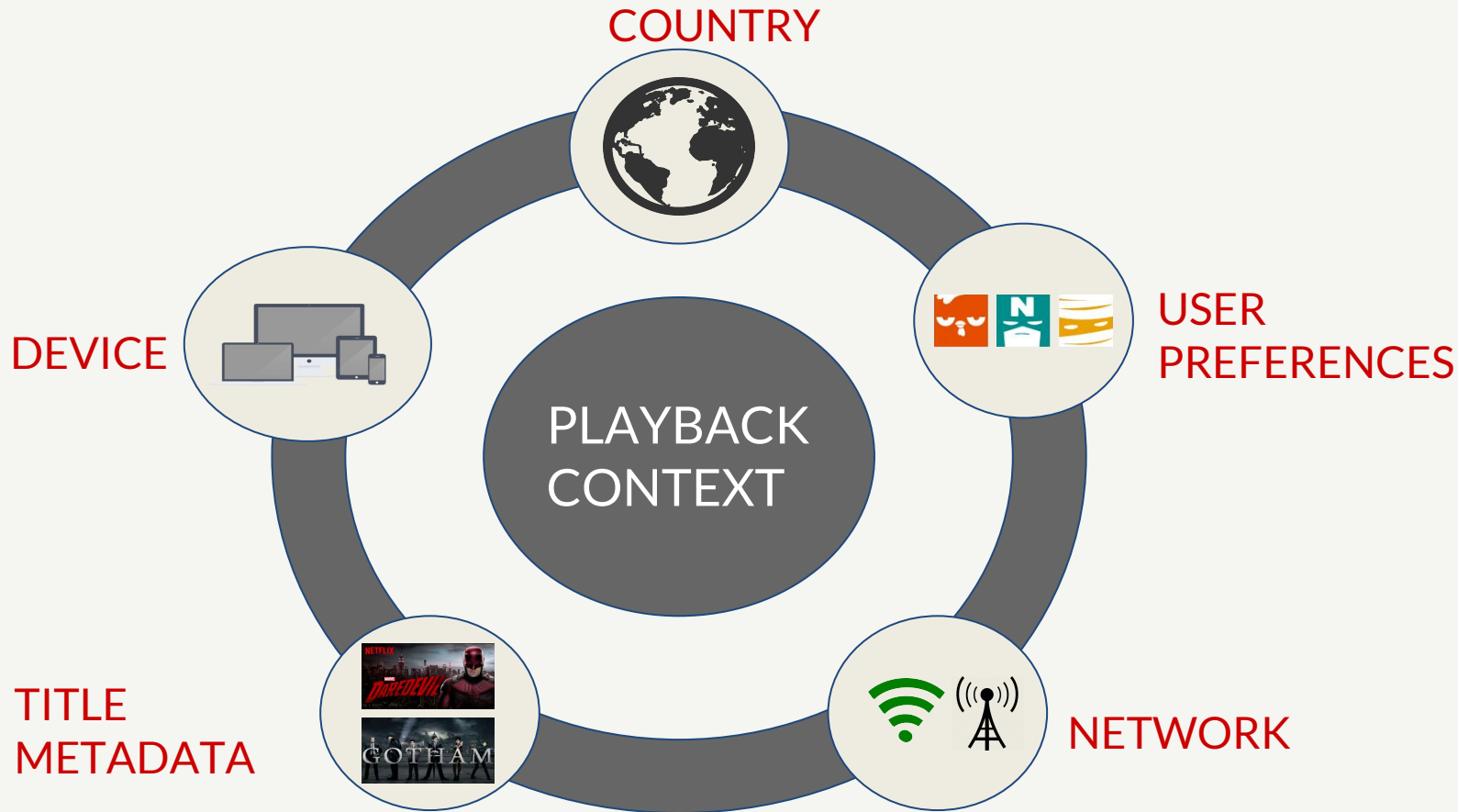Error Code: M7111-1331-2206

# Requirements - Playback Error

- Quick Incident Resolution
- Interactive Dashboards
- Near realtime Search
- Ad Hoc Queries

Interactive Exploration

# Top N queries

2+ Hours → Under 10 Minutes

# ★ Viewing History

# My Activity

See recent account access

| 11/7/17 | Stranger Things: Stranger Things 2: "Chapter Eight: The Mind Flayer" | Report a problem | ✕ |
| 11/3/17 | Stranger Things: Stranger Things 2: "Chapter Seven: The Lost Sister" | Report a problem | ✕ |
| 11/3/17 | Stranger Things: Stranger Things 2: "Chapter Nine: The Gate" | Report a problem | ✕ |
| 11/3/17 | Stranger Things: Stranger Things 2: "Chapter Six: The Spy" | Report a problem | ✕ |

# Requirements - Viewing History

- Time series dataset
- Support high writes
- Cross region replication
- Large dataset

# Growth of Viewing History

elastic

DYNOMITE

EVCache

cassandra

GUESS?

cassandra

➡ Multi-datacenter, multi-directional replication

➡ Highly availability and scalability

# Data Model

# New Data Model

# Requirements - DAM

- One backend plane for all asset metadata

- Storage of relationships/connected data

- Searchable

elastic

DYN*MITE

EVCache

cassandra

GUESS?

➡ Distributed GraphDB

➡ Support for various
storage backends

# ★Distributed Delayed Queues

# Requirements - Delayed Queues

- Distributed

- Highly concurrent

- At-least-once delivery semantics

- Delayed queue

- Priorities within the shard

# Data Model

For each queue three set of Redis data structures are maintained:

1. A Sorted Set containing queued elements by score.

2. A Hash set that contains message payload, with key as message ID.

3. A Sorted Set containing messages consumed by client but yet to be acknowledged. Un-ack set.

© DESPAIR.COM

CHALLENGES

I EXPECTED TIMES LIKE THIS - BUT I NEVER THOUGHT
THEY'D BE SO BAD, SO LONG, AND SO FREQUENT.

MAINTENANCE

# Current Approach

# Subject Matter Expert

# CDE Service

*"Empowering CDE to provide datastores as a service"*

# CDE Service

- Thresholds/SLAs
- Cluster metadata
- Self Service
- Contact information
- Maintenance windows

# Architecture

# SLA

**Cassandra »** cass_xyz

Customer View

PROD

## Customer Details

Customer Emails

abc@netflix.com

Customer Slack

cde

Customer PagerDuty Service

abc@netflix.com

DBEng Owners

Update Details

⌃ General Settings

⌃ Maintenance Windows

⌄ SLAs

| | | | |
|---|---|---|---|
| Read Latency (ms) | 99th 500 | 95th 200 | Avg 200 |
| Write Latency (ms) | 99th 500 | 95th 200 | Avg 200 |
| Disk Usage (%) | Fatal 80 | Warn 60 | |
| Co-Ordinator | # Reads 10000 | # Writes 10000 | |
| Node | # Reads 5000 | # Writes 5000 | |
| Max Row (bytes) | Size 10000 | | |

Update SLAs

⌃ Environment Attributes

# Cassandra Clusters  [10]  PROD: [6]  TEST: [4]

[Add new Cluster...]  [Edit Cluster Defaults]

Show [25] entries          Search: [____]   [Copy table to clipboard]  [Export to Excel]  [Show/Hide columns]

| Env | Region | Type | # Nodes | Customer Email | C* Version | C* JDK | Priam Version | Instance Type | Avg Node Size | Oldest Instance | EC2 Cost | S3 Primary Cost | S3 Secondary Cost |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **cas_xyz** ⁞ | | | | | | | | | | | 💲 EC2: [343] | S3: $ [34] | 📈 🖌 🔄 |
| prod | eu-west-1 | MR | 96 | abc@netflix.com | 2.1.17.1428… | JDK 8.0_45… | 6.84.0-h11… | i2.4xlarge(96) | 454.0 GB | 394 days | [2] | [4] | [34] |
| prod | us-east-1 | MR | 96 | abc@netflix.com | 2.1.17.1428… | JDK 8.0_45… | 6.84.0-h11… | i2.4xlarge(96) | 529.2 GB | 382 days | [2] | [3] | [34] |
| prod | us-west-2 | MR | 96 | abc@netflix.com | 2.1.17.1428… | JDK 8.0_45… | 6.84.0-h11… | i2.4xlarge(96) | 564.3 GB | 388 days | [5] | [6] | [65] |
| **cass_abc** ⇣ | | | | xyz@netflix.com | | | | | | | 💲 EC2: $ [ ] | S3: $ [ ] | 📈 🖌 🔄 |
| test | eu-west-1 | MR | 96 | | 2.1.17.1428… | NA(96) | 6.85.0-h11… | i2.2xlarge(96) | 487.0 GB | 6 days | [5] | [34] | [55] |
| test | us-east-1 | MR | 96 | xyz@netflix.com | 2.1.17.1428… | NA(96) | 6.85.0-h11… | i2.2xlarge(96) | 487.8 GB | 6 days | [65] | [56] | [56] |
| **cass_test** ⇣ | | | | | | | | | | | 💲 EC2: $ [ ] | | 📈 🖌 🔄 |
| test | us-east-1 | Island | 6 | | 2.1.17.1428… | NA(6) | 6.84.0-h11… | i2.xlarge(6) | 1.8 GB | 450 days | [ ] | | |

# Create a new Elasticsearch Cluster

Create Cluster

**Before you begin:**

- Elasticsearch is not recommended as a primary data store; if you choose to use it as one, please make sure to take steps to prevent data loss.
- If your use case is not large enough for a dedicated cluster, consider creating your index on our shared cluster "es_share5" instead.
- For more information about working with Elasticsearch, please see http://go/elasticsearch

**Cluster Name**

es_

**Owners**

Select a mailgroup or user email    ⌄

**Cluster Topology** ⑦

◉ Island Clusters    ◯ Tribe Configuration

**PROD**

Regions to create **Island Clusters** in

☐ eu-west-1    ☐ us-east-1    ☐ us-west-2

Estimated data size per region (GB) ⑦

Size in GB

**+ Show Advanced Options**

**TEST**

Regions to create **Island Clusters** in

☐ eu-west-1    ☐ us-east-1

# Dynomite Goals

Edit

## Dynomite v0.5.9

| ALL | PROD | TEST |
|-----|------|------|
| **98%** | **100%** | **92%** |

## Dynomite v0.6

| ALL | PROD | TEST |
|-----|------|------|
| **2%** | **0%** | **6%** |

## Xenial Upgrade

| ALL | PROD | TEST |
|-----|------|------|
| **100%** | **100%** | **100%** |

# Machine learning

# Pattern in Disk usage

Rolling Mean & Standard Deviation

```
Results of Dickey-Fuller Test:
Test Statistic                    0.350320
p-value                           0.979539
#Lags Used                       10.000000
Number of Observations Used     248.000000
Critical Value (5%)              -2.873266
Critical Value (1%)              -3.456996
Critical Value (10%)             -2.573019
dtype: float64
```

# Cde Channel

2 members | Add a topic

Today

36.1503231023 in 90 days

`cass_xyz` is `18` nodes with current read latency of `17203.4597222` and may miss read latency with expected value of `17267.4362002` in 90 days

`cass_xyz` is `18` nodes with current disk usage of `14.1122366141` and may reach disk usage of `42.0660391165` in 90 days

`cass_xyz` is `18` nodes with current read latency of `14673.5980873` and may miss read latency with expected value of `14745.335325` in 90 days

`cass_xyz` is `24` nodes with current disk usage of `16.9526726339` and may reach disk usage of `39.7745393664` in 90 days

`cass_xyz` is `12` nodes with current disk usage of `5.15693085154` and may reach disk usage of `36.7731174652` in 90 days

`cass_xyz` is `12` nodes with current disk usage of `5.18868543363` and may reach disk usage of `37.165458283` in 90 days

`cass_xyz` is `12` nodes with current disk usage of `29.87` and may reach disk usage of `74.99` in 90 days

`cass_xyz` is `23` nodes with current disk usage of `11.49` and may reach disk usage of `34.17` in 90 days

`cass_xyz` is `12` nodes with current disk usage of `29.89` and may reach disk usage of `73.60` in 90 days

`cass_xyz` is `36` nodes with current disk usage of `19.68` and may reach disk usage of `36.13` in 90 days

`cass_xyz` is `18` nodes with current read latency of `17203.46` and may miss read latency with expected value of `17267.44` in 90 days

`cass_xyz` is `18` nodes with current disk usage of `14.11` and may reach disk usage of `42.07` in 90 days
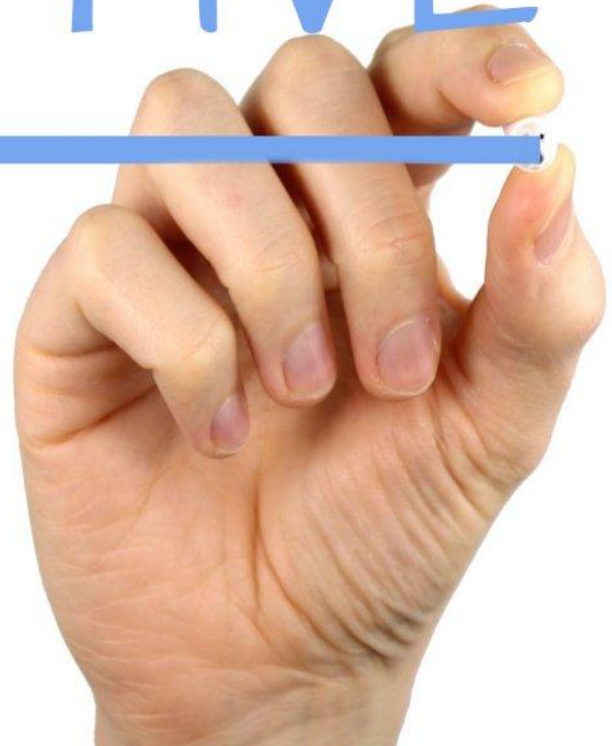
`cass_xyz` is `24` nodes with current disk usage of `16.95` and may reach disk usage of `39.77` in 90 days

`cass_xyz` is `12` nodes with current disk usage of `5.16` and may reach disk usage of `36.77` in 90 days

`cass_xyz` is `12` nodes with current disk usage of `5.19` and may reach disk usage of `37.17` in 90 days

`cass_xyz` is `18` nodes with current disk usage of `16.38` and may reach disk usage of `45.66` in 90 days
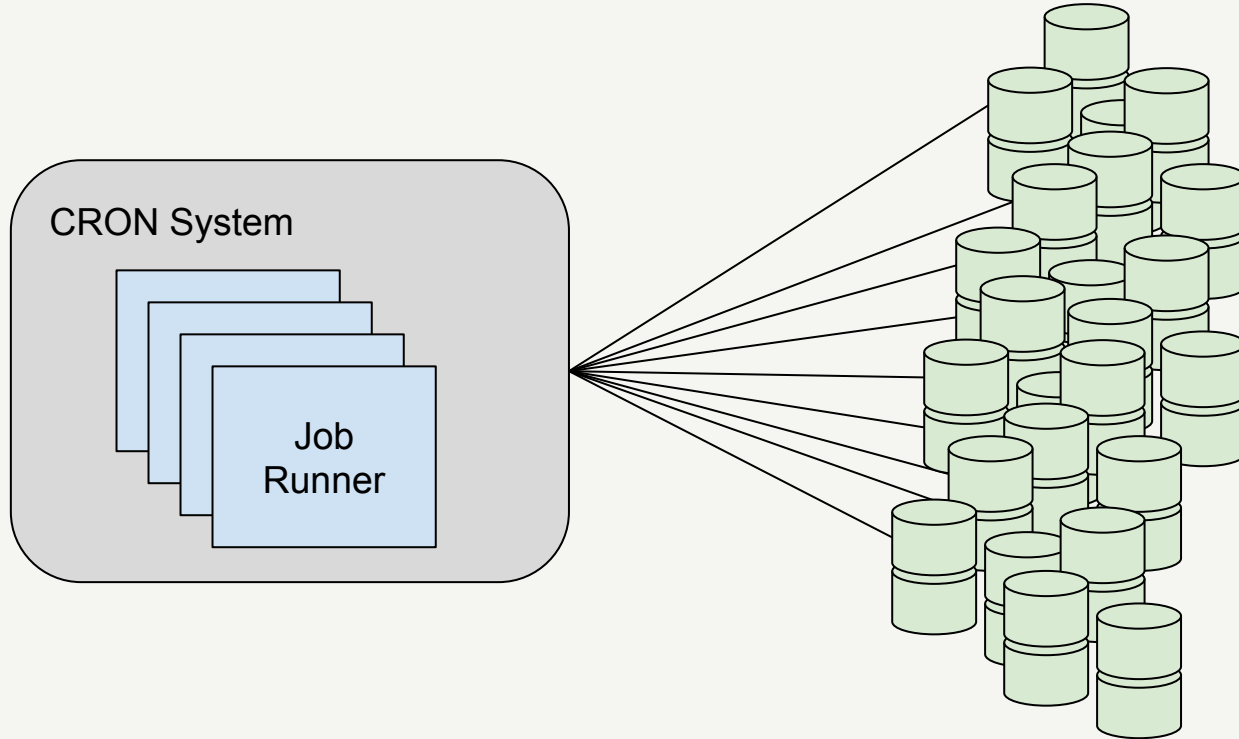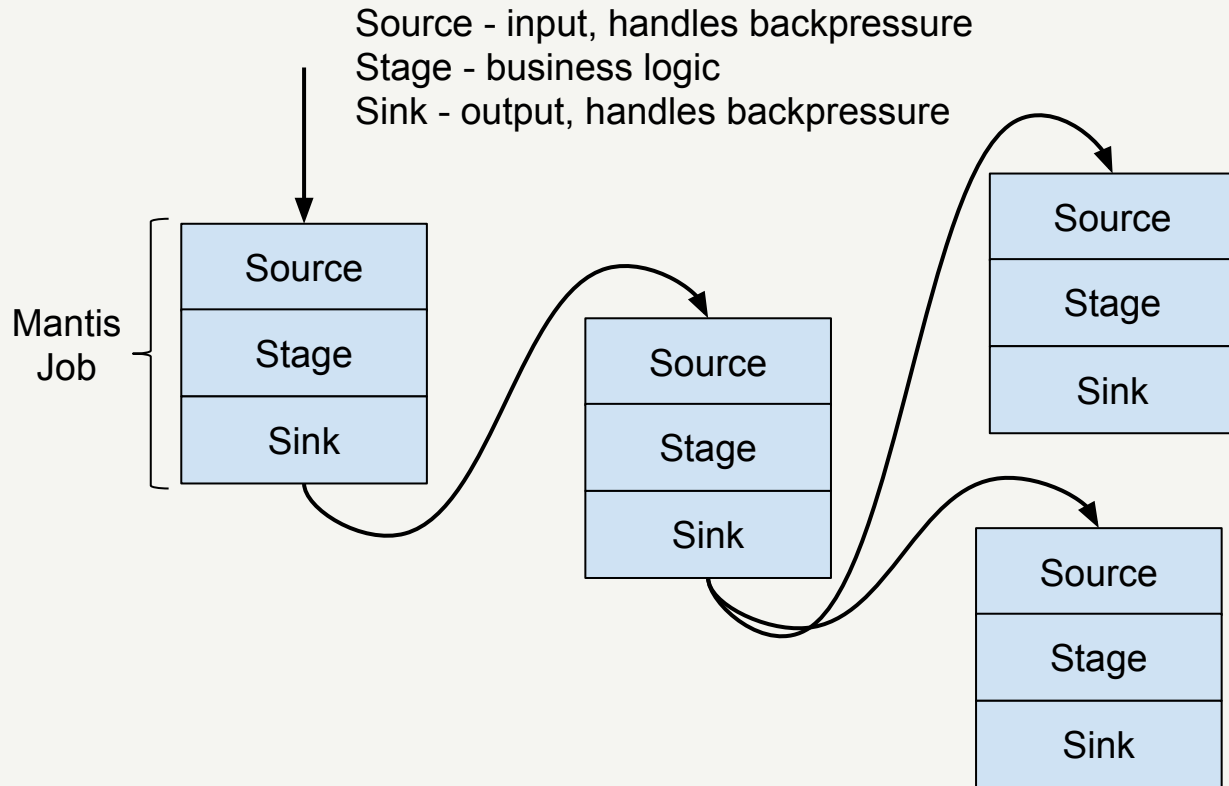
UPGRADE

# Common Approach

# Streaming micro-services

Source - input, handles backpressure
Stage - business logic
Sink - output, handles backpressure

Mantis Job

| Source |
| --- |
| Stage |
| Sink |

| Source |
| --- |
| Stage |
| Sink |

| Source |
| --- |
| Stage |
| Sink |

| Source |
| --- |
| Stage |
| Sink |

# Real Time Dash (Cluster View)

# Takeaway

# Talk: Microservices: Patterns and Practices Panel

**A** Track: **Microservices: Patterns and Practices**

**♀** Location: **Ballroom A**

**⊙** Day of week: **Tuesday**

**⊙** Duration: **4:10pm - 5:00pm**

Microservices almost seem to be the de facto way to build systems today, but are they always the answer? If they are the answer, what are the challenges you'll face at scale (both from a technical and organizational level)? What are the strategies you should use now that you are effectively building a distributed system? ...or what's the one thing you wish you'd known before you got here? These questions and more will be asked in the Microservices: Pattern's & Practices Ask Me Anything or AMA (a significant portion of the time will be available for the audience to get their questions answered as well). This session joins together many of the conference's most popular sessions speakers with the trackhost from the Microservices track to have a frank and honest discussion on Microservices. Join us to have your Microservices questions answered.