# Executive Briefing: What You Need to Know about Fast Data

Dean Wampler, Ph.D.
dean@lightbend.com
@deanwampler
polyglotprogramming.com/talks

Lightbend

# Based on this report

go.lightbend.com/fast-data-architectures-for-streaming-applications-oreilly-2nd-edition
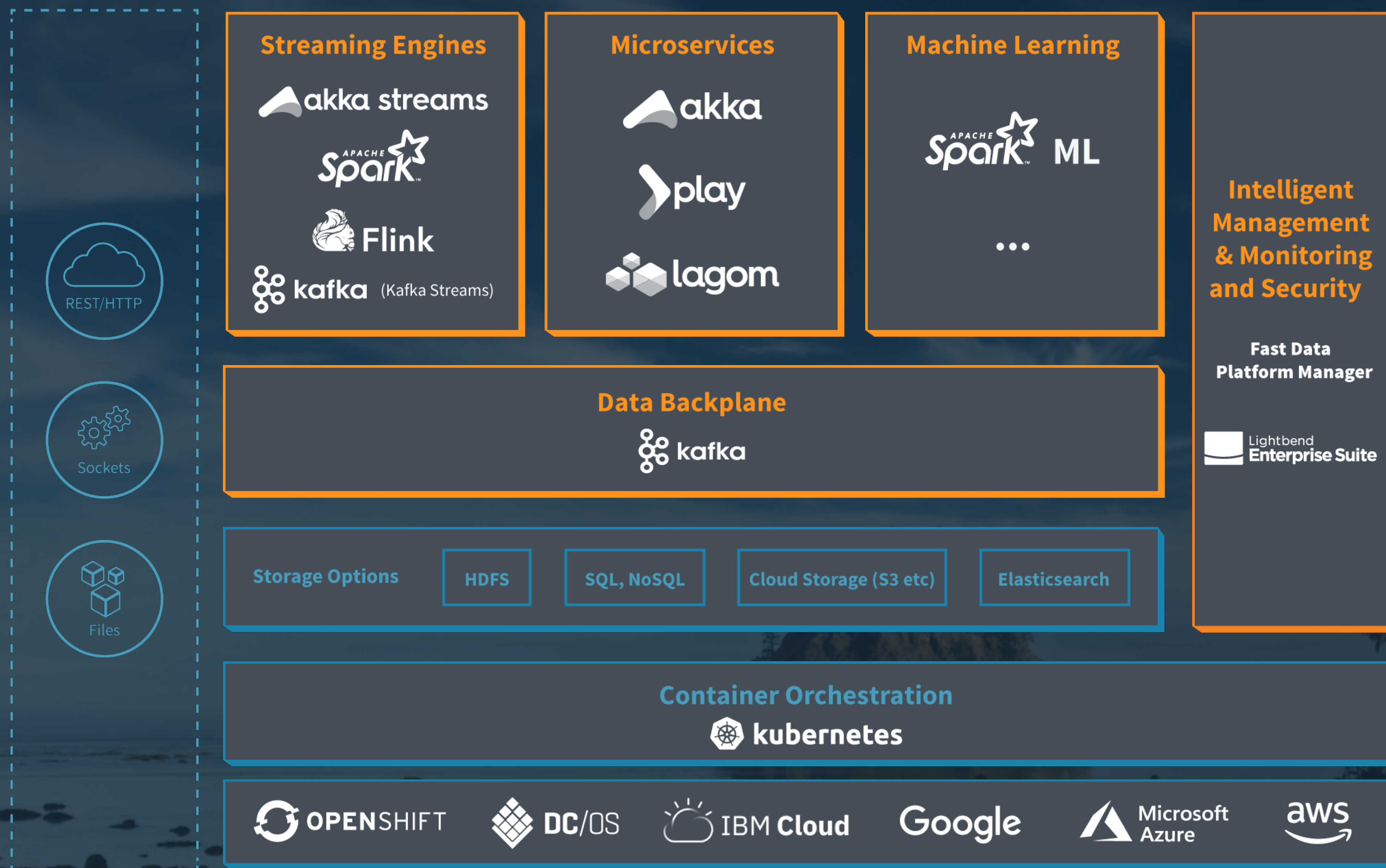
O'REILLY®

# Fast Data Architectures for Streaming Applications

## Getting Answers Now from Data Sets that Never End

Dean Wampler

I lead the Lightbend Fast Data Platform project; streaming data and microservices

lightbend.com/fast-data-platform

What We'll Discuss

- Why streaming? Why now?

- How to choose technologies

- The impact streaming will have on your organization

What We'll Discuss

Why Streaming?

- New opportunities that require streaming
  - Media content is obviously one ;)
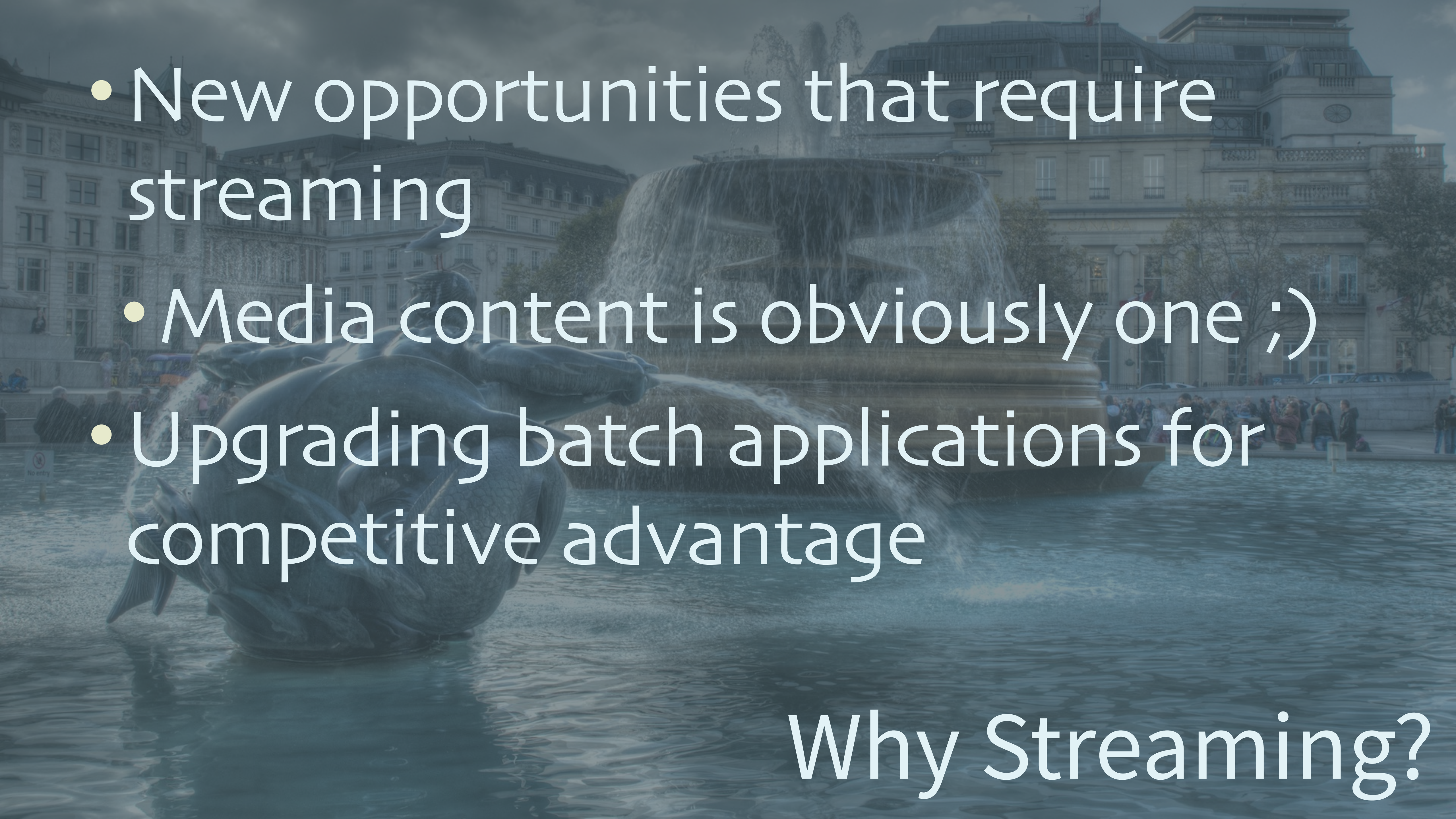  - Upgrading batch applications for competitive advantage

Why Streaming?

# Fast Data Use Cases

**Similar IoT Architectures**

| Predictive Analytics | IoT | Real-time Personalization | Real-time Financial Processes |
|---|---|---|---|
| Apply ML models to large volumes of device data to pre-empt failures / outages | Real-time consumer and industrial Device and Supply Chain management at scale | Real-time marketing based on behavior, location, inventory levels, product promotions, etc. | Drive better business outcomes through real-time risk, fraud detection, compliance, audit, governance, etc. |
| Hewlett Packard Enterprise | STARBUCKS COFFEE | Royal Caribbean INTERNATIONAL | Capital One |

https://www.lightbend.com/customers

# Predictive Analytics

**Hewlett Packard**
Enterprise

- ML models applied to device telemetry to detect anomalies

- Preemptive maintenance prevents potential failures that would impact users

# Example Architecture

Device Session Microservices

10   Corrective Action

Microservice
Microservice
**Microservice**

Device

1   Telemetry

9   Ingest Scores

Broker

Kafka Cluster

Persistence

6, 7, 8

6. → Data Pipeline
7. → Model Serving
8. ← Anomalies

Akka Streams

Kafka Streams

…

Low Latency Microservices

2, 3

2. → Model Training
3. ← New models

4. ← Model Storage
5. → Boot up, historical data

4, 5

Spark

Mini-batch, Batch

→ Sessions
→ Streams
→ Storage

Data Center

# Example Architecture

Device Session
Microservices

Session management,
REST microservices

Model
Scoring

10  Corrective
    Action
                        Microservice
                        Microservice
                        Microservice

                                              Akka Streams

                                              Kafka Streams

                                              ...

Device

9  Ingest Scores        6, 7, 8

                        6. → Data Pipeline       Low Latency
                        7. → Model Serving       Microservices
                        8. ← Anomalies

1  Telemetry

Broker                  2, 3                     Model
                                                 Training

                        2. → Model Training
Kafka Cluster           3. ← New models

                                                 Spark

                        4. ← Model Storage
                        5. → Boot up,
                        historical data          Mini-batch, Batch

Sessions

Streams

Storage

Persistence             4, 5

                        Three groups of functionality

# Example Architecture



**Device Session Microservices**

Microservice
Microservice
Microservice

10 Corrective Action

Device

1 Telemetry

9 Ingest Scores

Broker

Kafka Cluster

Ingest device telemetry to Kafka

→ Storage

Persistence

6, 7, 8
6. → Data Pipeline
7. → Model Serving
8. ← Anomalies

Akka Streams
Kafka Streams
...

Low Latency Microservices

2, 3
2. → Model Training
3. ← New models

4. ← Model Storage
5. → Boot up, historical data

4, 5

Spark

Mini-batch, Batch

Data Center

# Example Architecture

Device Session Microservices

Microservice
Microservice
**Microservice**

⑩ Corrective Action

Akka Streams

⑨ Ingest Scores

⑥, ⑦, ⑧

Read telemetry into a periodic Spark job for model training to detect anomalies

**Device**

① Telemetry

6. → Data Pipeline
7. → Model Serving
8. ← Anomalies

**Broker**

Kafka Cluster

②, ③

2. → Model Training
3. ← New models

**Spark**

Mini-batch, Batch

→ Sessions
→ Streams
→ Storage

4. ← Model Storage
5. → Boot up, historical data

④, ⑤

Persistence

Large data volume, Long latency (seconds-days)

# Example Architecture

# Example Architecture

Device Session Microservices

**10** Corrective Action

Microservice
Microservice
**Microservice**

**9** Ingest Scores

Akka Streams
Kafka Streams
...

Low Latency Microservices

**6, 7, 8**
6. → Data Pipeline
7. → Model Serving
8. ← Anomalies

Device

**1** Telemetry

Broker

Kafka Cluster

**2, 3**
2. → Model Training
3. ← New models

→ Sessions
→ Streams
→ Storage

4. ← Model Storage
5. → Boot up, historical data

**4, 5**

Persistence

Updated model parameters are also written to longer-term storage (for system restarts, auditing, …)

# Example Architecture

# Example Architecture

Session management, REST microservices

Model Scoring

Model Training

Three groups of functionality

Device Session Microservices

Microservice
Microservice
Microservice

Corrective Action

10

Device

1 Telemetry

9 Ingest Scores

Broker

Kafka Cluster

Persistence

Akka Streams
Kafka Streams
...

Low Latency Microservices

6, 7, 8

6. → Data Pipeline
7. → Model Serving
8. ← Anomalies

2, 3

2. → Model Training
3. ← New models

4. ← Model Storage
5. → Boot up, historical data

4, 5

Spark

Mini-batch, Batch

Sessions
Streams
Storage

- Integration of Machine Learning/Artificial Intelligence with streaming is a common challenge right now

Device Session Microservices

Microservice

Microservice

Microservice

Corrective Action

10

Device

1

Telemetry

Ingest Scores

9

Broker

Kafka Cluster

Persistence

Sessions

Streams

Storage

Akka Streams

Kafka Streams

...

Low Latency Microservices

6, 7, 8

6. → Data Pipeline
7. → Model Serving
8. → Anomalies

2, 3

2. → Model Training
3. ← New models

4. ← Model Storage
5. → Boot up, historical data

4, 5

Spark

Mini-batch, Batch

Data Center

- Network overhead for telemetry ingestion too high?

- Model serving latency too long?

- Datacenter unavailable?

- Idea: Serve models on the device!

# Internet of Things

- Real-time consumer and industrial device and supply chain management at scale

Example Architecture

Device Session Microservices

Microservice
Microservice
Microservice

10 Corrective Action

Device

9 Ingest Scores

1 Telemetry

Broker

Kafka Cluster

Akka Streams
Kafka Streams
...

Low Latency Microservices

6, 7, 8
6. → Data Pipeline
7. → Model Serving
8. ← Anomalies

2, 3
2. → Model Training
3. ← New models

4. ← Model Storage
5. → Boot up, historical data

4, 5

Spark

Mini-batch, Batch

Persistence

Sessions
Streams
Storage

What we just discussed...

Data Center

Edge-Scoring Example Architecture

Edge-Scoring Example Architecture

Edge-Scoring Example Architecture

Edge-Scoring Example Architecture

Edge-Scoring Example Architecture

Edge-Scoring Example Architecture

# Edge-Scoring Example Architecture

**Device Session Microservices**

Microservice

Microservice

Microservice

8 Corrective Action

Device

7 Model Parameters

1 Telemetry

6 Model Parameters

Broker

Kafka Cluster

2, 3

2. → Model Training
3. ← New models

Spark

Mini-batch, Batch

4. ← Model Storage
5. → Boot up, historical data

4, 5

Persistence

→ Sessions
→ Streams
→ Storage

Recap: Edge Serving

Technology Choices

- More than "faster" Hadoop…

- New architectures that merge data processing with microservices

Technology Choices

Recall Hadoop…

- Data warehouse replacement

- Historical analysis

- Interactive exploration

- Offline training of machine learning models

- …

**Compute**

- MapReduce jobs
- Spark jobs
- ...

**Resource Management**

submit to...

**YARN**

| Master Node | Worker Node #1 | #2 | ... |
| Resource Manager | Node Manager | | |

**HDFS**

| Name Node | Data Node | | |

Disk

**Storage**

Optimized for storing lots of data *at rest*, with subsequent processing, but not optimized for data *in motion*.

- Hadoop is ideal for batch and interactive apps

- ... but also constrained by that model

New Fast Data
Architecture

Kubernetes, Mesos, YARN, ...
Cloud or on-premise

Reactive Platform
Go  Node.js  ...
Microservices

ZK
ZooKeeper Cluster

REST

Sockets

Files

Broker
Kafka Cluster

Spark
Flink
Beam

Akka Streams
Kafka Streams
Low Latency

Disk
HDFS

S3, ...

Spark
Beam
Mini-batch

Search

SQL/
NoSQL
Persistence

Spark
...
Batch

Events
Streams
Storage

Flesh out earlier
example architectures

Reactive Platform
Go | Node.js | ...
Microservices

③

ZK
ZooKeeper Cluster

④

Kubernetes, Mesos, YARN, ...
Cloud or on-premise

"Events" - e.g., REST messages, sessions, alerts, ...

REST

②

①

Sockets

Files

Broker
Kafka Cluster

⑤

⑥

⑦

⑧

A
Ka

"Streams" - one-way data flows, e.g., sockets or files, including logs, metrics, other telemetry, click streams, etc.

Events

Streams

Storage

Disk
HDFS

S3, ...

Search

SQL/NoSQL

Persistence

⑨

Spark
...
Batch

Kubernetes, Mesos, YARN, ...
Cloud or on-premise

Reactive Platform
Go    Node.js    ...
Microservices

③

ZK
ZooKeeper Cluster

④

"Events" - e.g., REST messages, sessions, alerts, …

REST

②

⑥

Sockets

①

Broker
Kafka Cluster

A
Ka

"Streams" - one-way data flows, e.g., sockets or files, including logs, metrics, other telemetry, click streams, etc.

Files

⑦

⑤

⑧

D  D  Disk

S3, ...

HDFS

Events

Streams

Storage

⑨

Search

SQL/

"Storage" - JDBC, async reads/writes to storage

Batch

Each has different volumes, velocities, latency characteristics, protocols, etc.

Lots of streaming engine options... why so many?

# You need choices

- Latency: how low?

- Volume per unit time: how high?

- Data processing: which kinds?

- Build, deploy, and manage services: what are your preferences?

# Why Microservices in Fast Data?

1. The trend is to run everything in big clusters using Kubernetes or Mesos

   - In the cloud or on-premise

# Why Microservices in Fast Data?
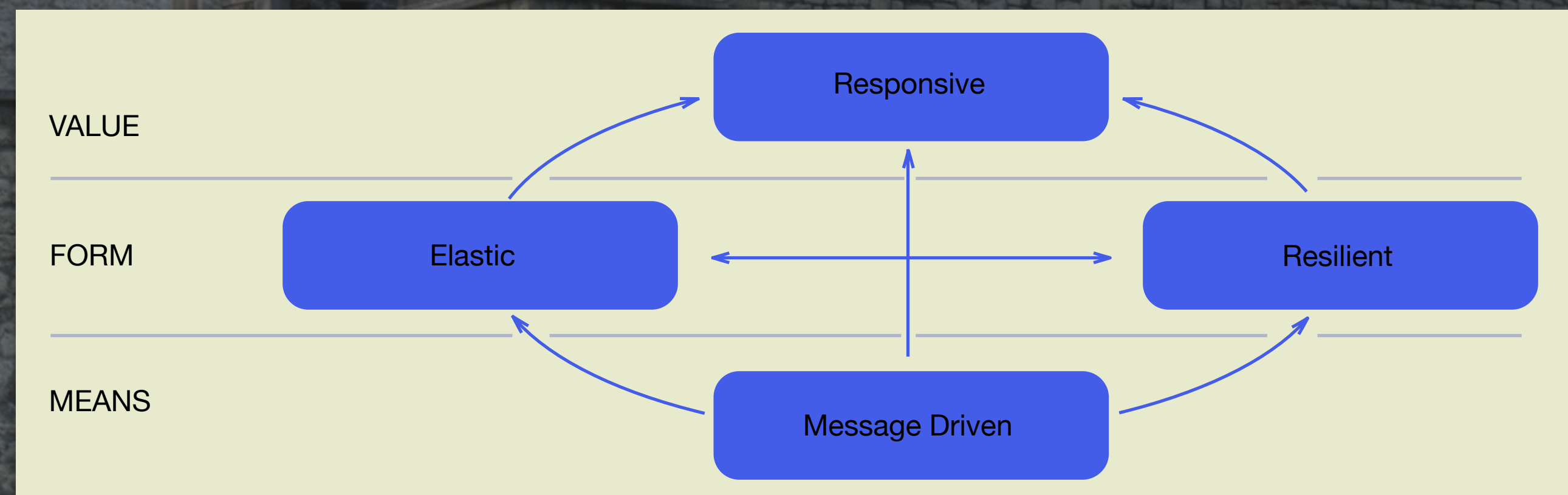
2. If streaming gives you information faster…

- … you'll want quick access to it in your other services!

# Why Microservices in Fast Data?

3. Streaming raises the bar on data services

- Compared to batch services, long-running streaming services must be more:

  - Scalable

  - Resilient

  - Flexible



| | | |
|---|---|---|
| VALUE | | Responsive |
| FORM | Elastic | Resilient |
| MEANS | | Message Driven |

https://www.reactivemanifesto.org/

# Why Microservices in Fast Data?

4. This leads to our last major point…

# Organizational Impact

# Organizational Impact

- Data engineers have to become good at highly-available microservices

- Microservice engineers have to become good at data

- … and Data scientists have to understand production issues

# The Past

Services     Big Data

Some overlap in concerns, architecture

# The Present

## Microservices & Fast Data

Much more overlap

# Lightbend
# Fast Data Platform

lightbend.com/fast-data-platform

**Streaming Engines**
akka streams
Apache Spark
Flink
kafka (Kafka Streams)

**Microservices**
akka
play
lagom

**Machine Learning**
Apache Spark ML
...

**Intelligent Management & Monitoring and Security**
Fast Data Platform Manager
Lightbend Enterprise Suite

REST/HTTP
Sockets
Files

**Data Backplane**
kafka

Storage Options | HDFS | SQL, NoSQL | Cloud Storage (S3 etc) | Elasticsearch

**Container Orchestration**
kubernetes

OPENSHIFT | DC/OS | IBM Cloud | Google | Microsoft Azure

What we discussed

Plus management and monitoring tools

lightbend.com/fast-data-platform

# lightbend.com/fast-data-platform

Dean Wampler, Ph.D.
dean@lightbend.com
@deanwampler
polyglotprogramming.com/talks

Lightbend