

Executive Briefing: What You Need to Know about Fast Data

Dean Wampler, Ph.D.
dean@lightbend.com
[@deanwampler](https://twitter.com/deanwampler)
polyglotprogramming.com/talks





Based on
this report

lightbend.com/fast-data-platform

O'REILLY®

Fast Data Architectures for Streaming Applications

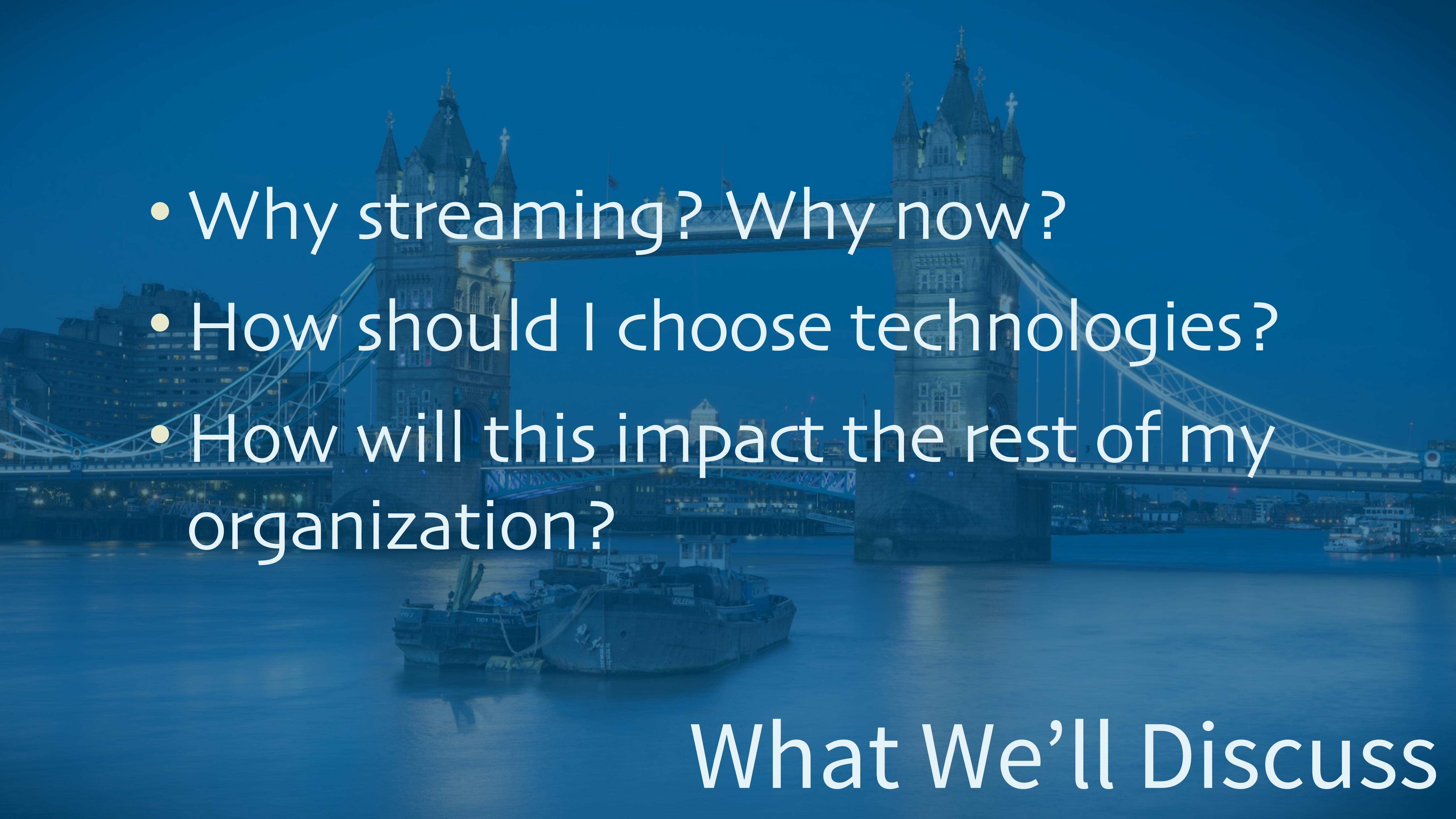
Getting Answers Now from
Data Sets that Never End

Dean Wampler

Compliments of
 Lightbend



What We'll Discuss

- 
- Why streaming? Why now?
 - How should I choose technologies?
 - How will this impact the rest of my organization?

What We'll Discuss



Why Streaming?

- 
- A large, stylized blue whale sculpture is the central focus of a city fountain. The whale is depicted in a dynamic, leaping pose, with its body curved and its tail pointing upwards. Water is spraying from its blowhole and along its back. In the background, there are classical buildings, including one with a prominent clock tower and another with a sign that reads "CANADIAN PACIFIC". People are visible around the fountain and on the surrounding stone walls.
- New opportunities that require streaming
 - Upgrading batch applications for competitive advantage

Why Streaming?



Fast Data Use Cases

Predictive Analytics

Apply ML models to large volumes of device data to pre-empt failures / outages



IoT

Real-time consumer and industrial Device and Supply Chain management at scale



Real-time Personalization

Real-time marketing based on behavior, location, inventory levels, product promotions, etc.



Real-time Financial Processes

Drive better business outcomes through real-time risk, fraud detection, compliance, audit, governance, etc.



Legacy Modernization

Accelerate decision making processes and optimize infrastructure costs by moving from batch to streaming



More at: <https://www.lightbend.com/customers>

Similar IoT Architectures

Fast Data Use Cases

Predictive Analytics

Apply ML models to large volumes of device data to pre-empt failures / outages

 Hewlett Packard Enterprise

IoT

Real-time consumer and industrial Device and Supply Chain management at scale



Real-time Personalization

Real-time marketing based on behavior, location, inventory levels, product promotions, etc.



Real-time Financial Processes

Drive better business outcomes through real-time risk, fraud detection, compliance, audit, governance, etc.

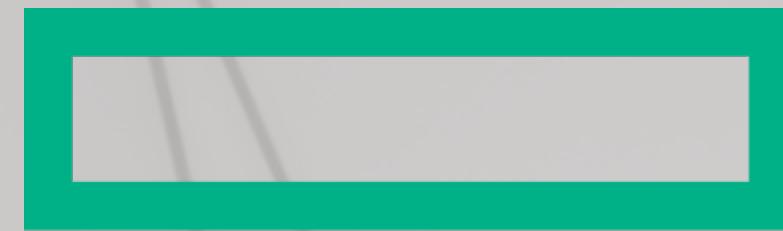


Legacy Modernization

Accelerate decision making processes and optimize infrastructure costs by moving from batch to streaming



More at: <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

Core Idea

Train models to look for anomalies... and score incoming telemetry.

Anomaly Handler

Corrective Actions

Probable Anomalies

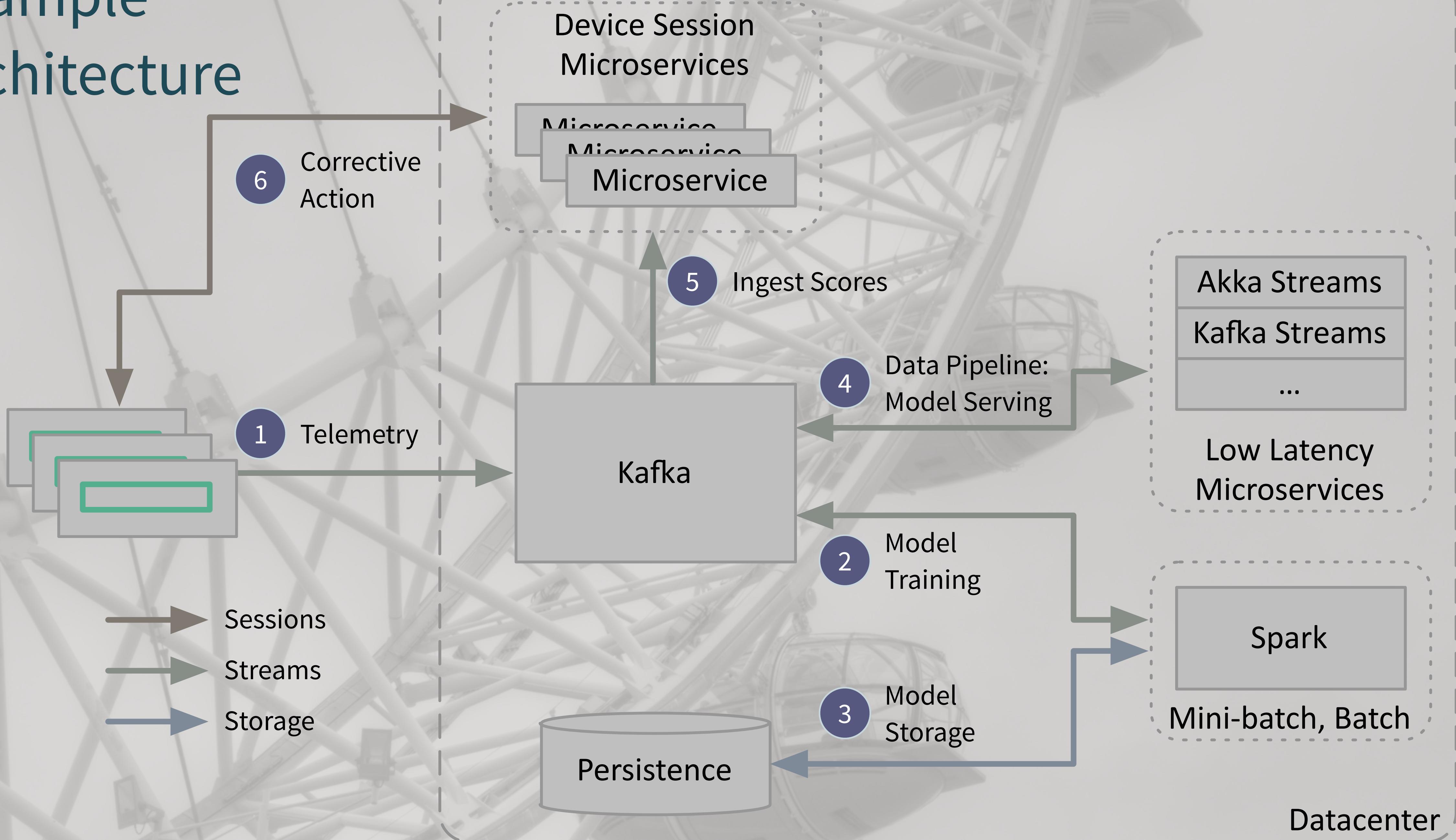
Anomaly Detection: Model

Ingest telemetry from edge devices.

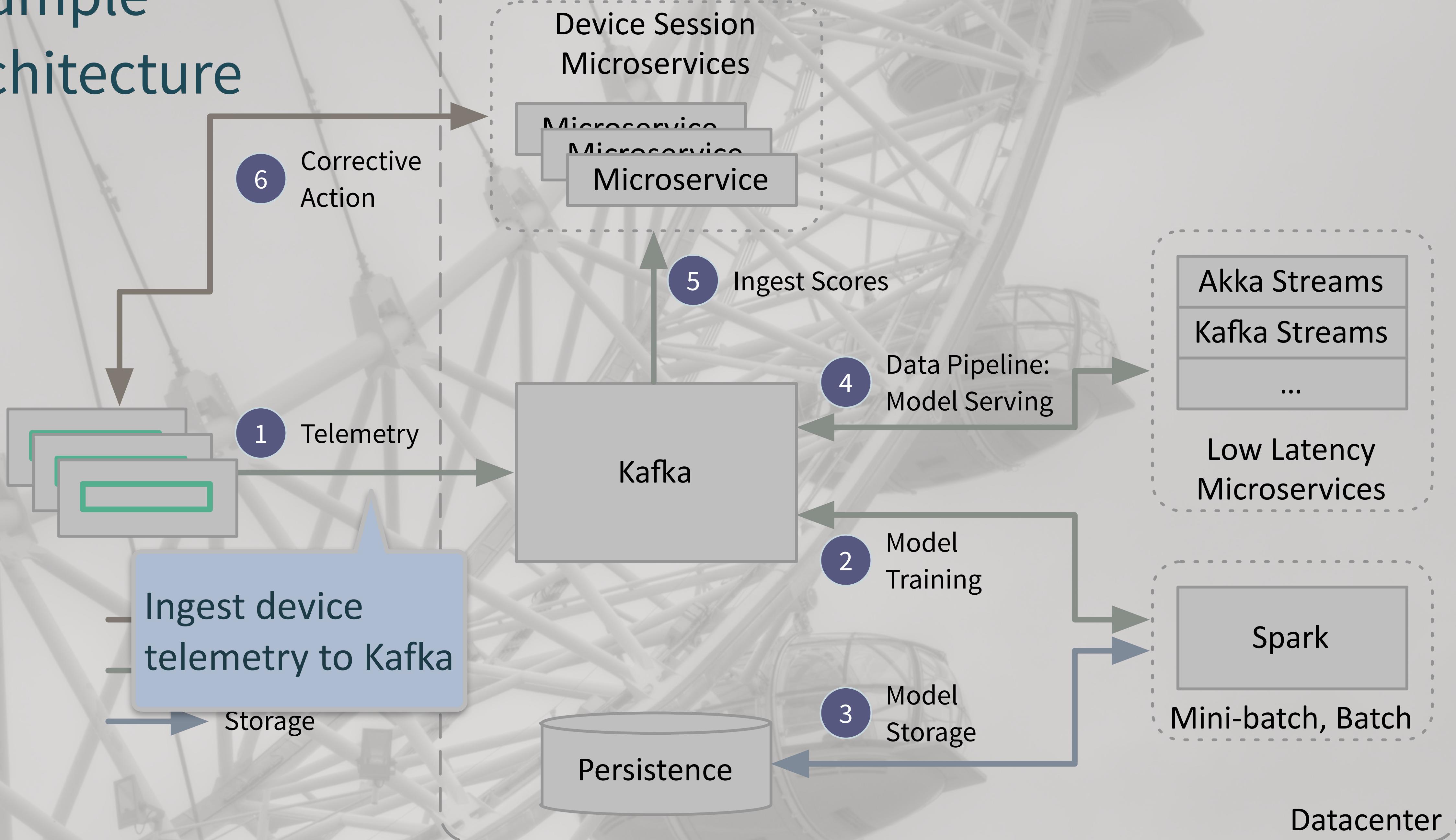
Nimble Storage



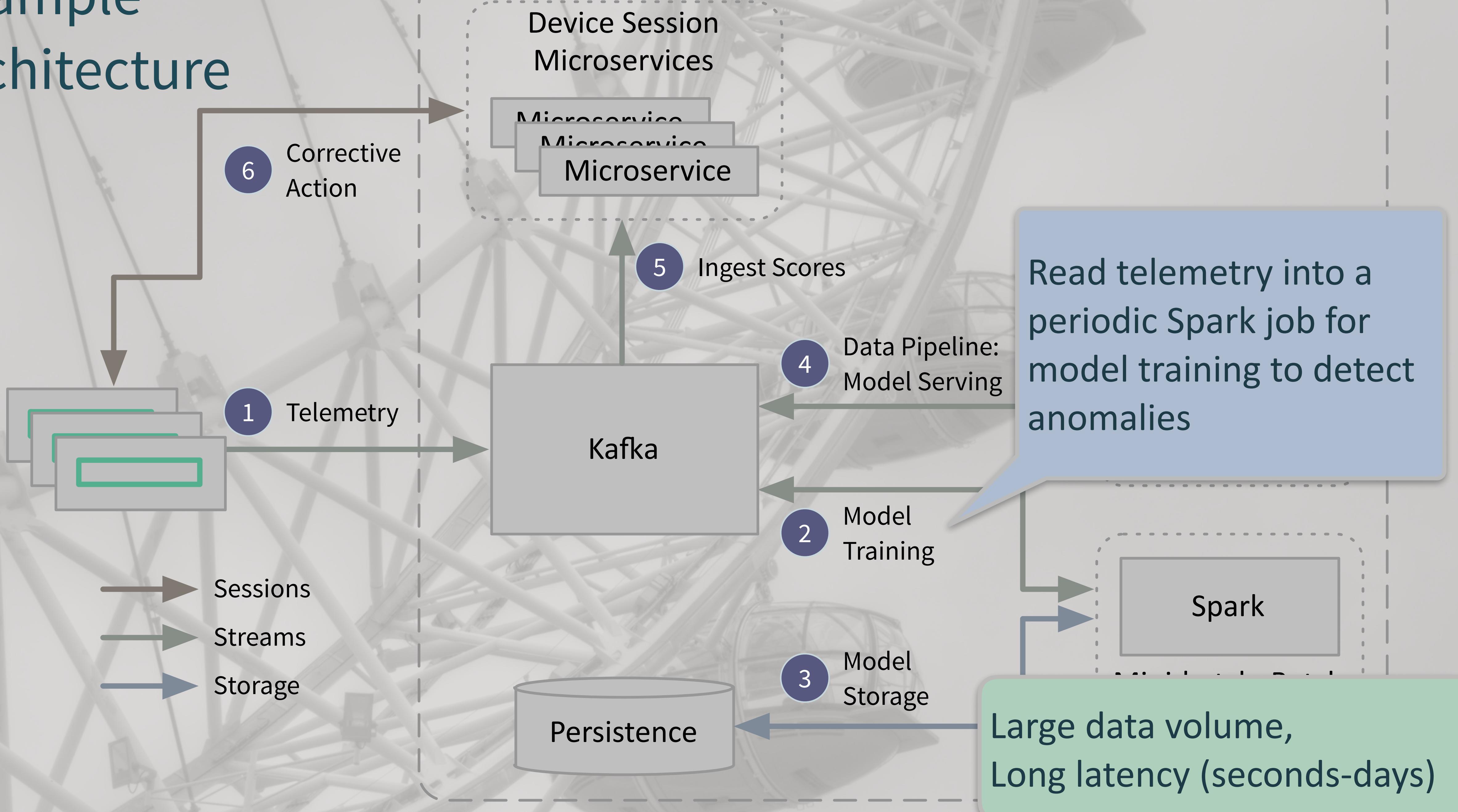
Example Architecture



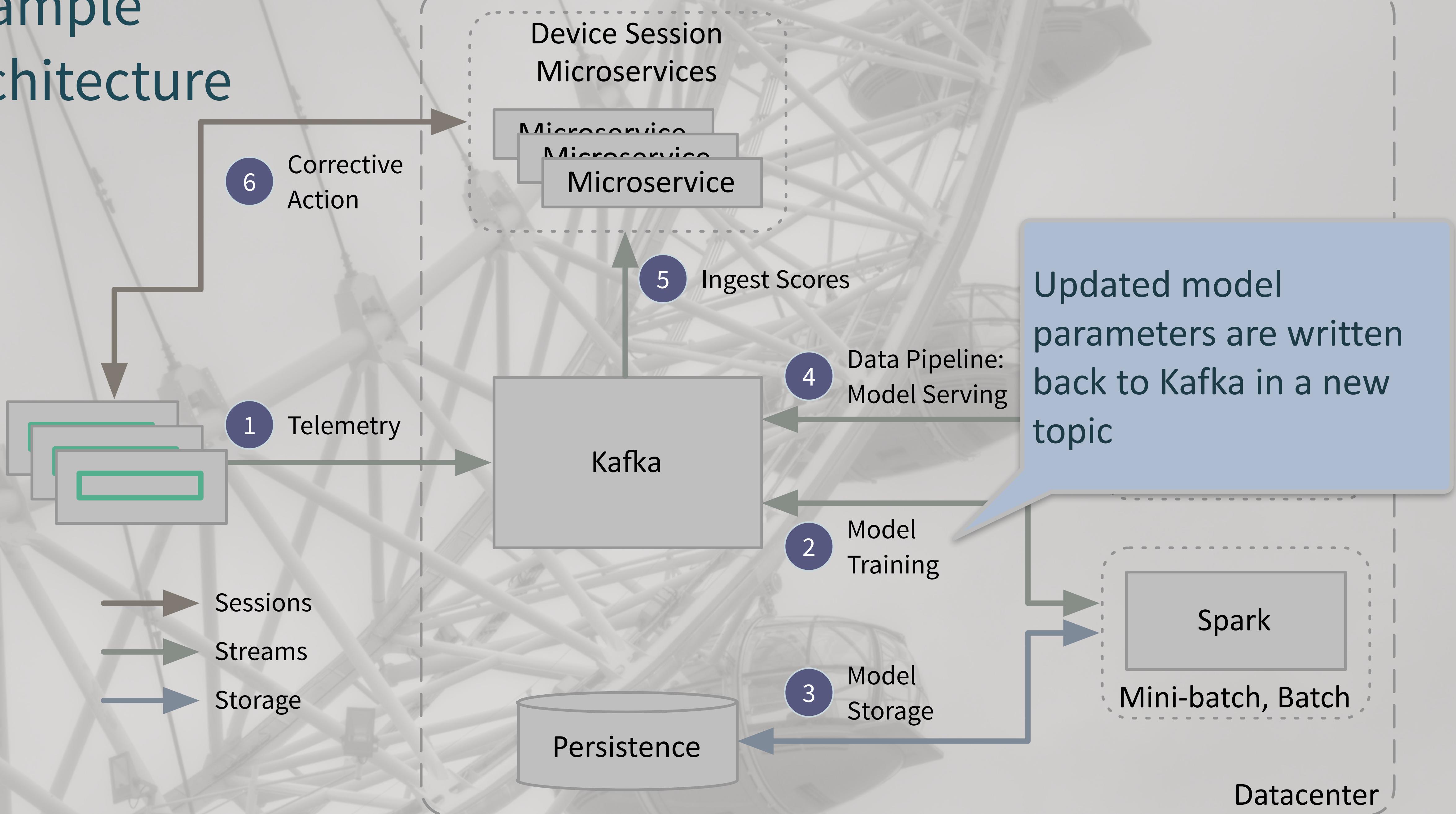
Example Architecture



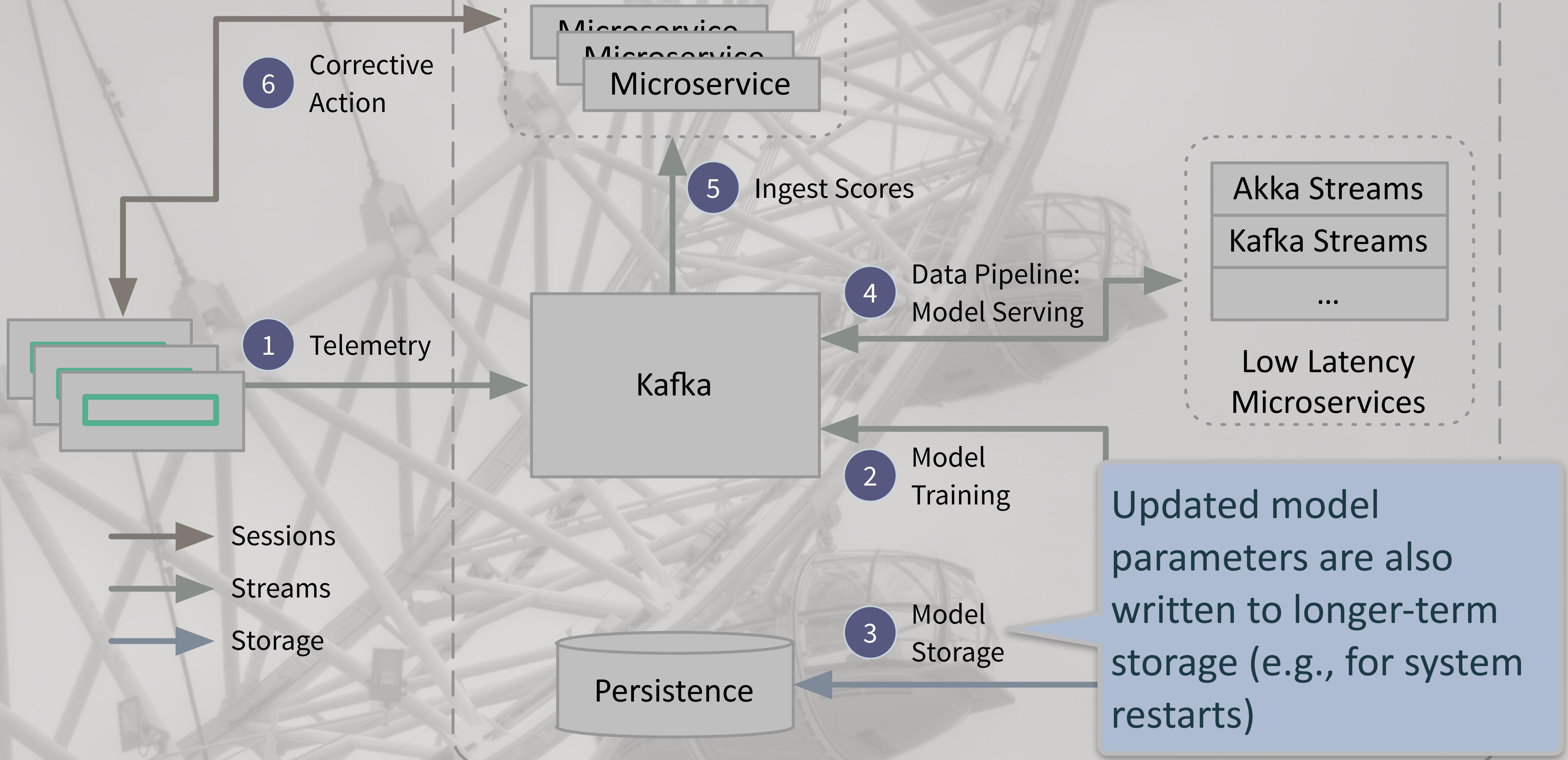
Example Architecture



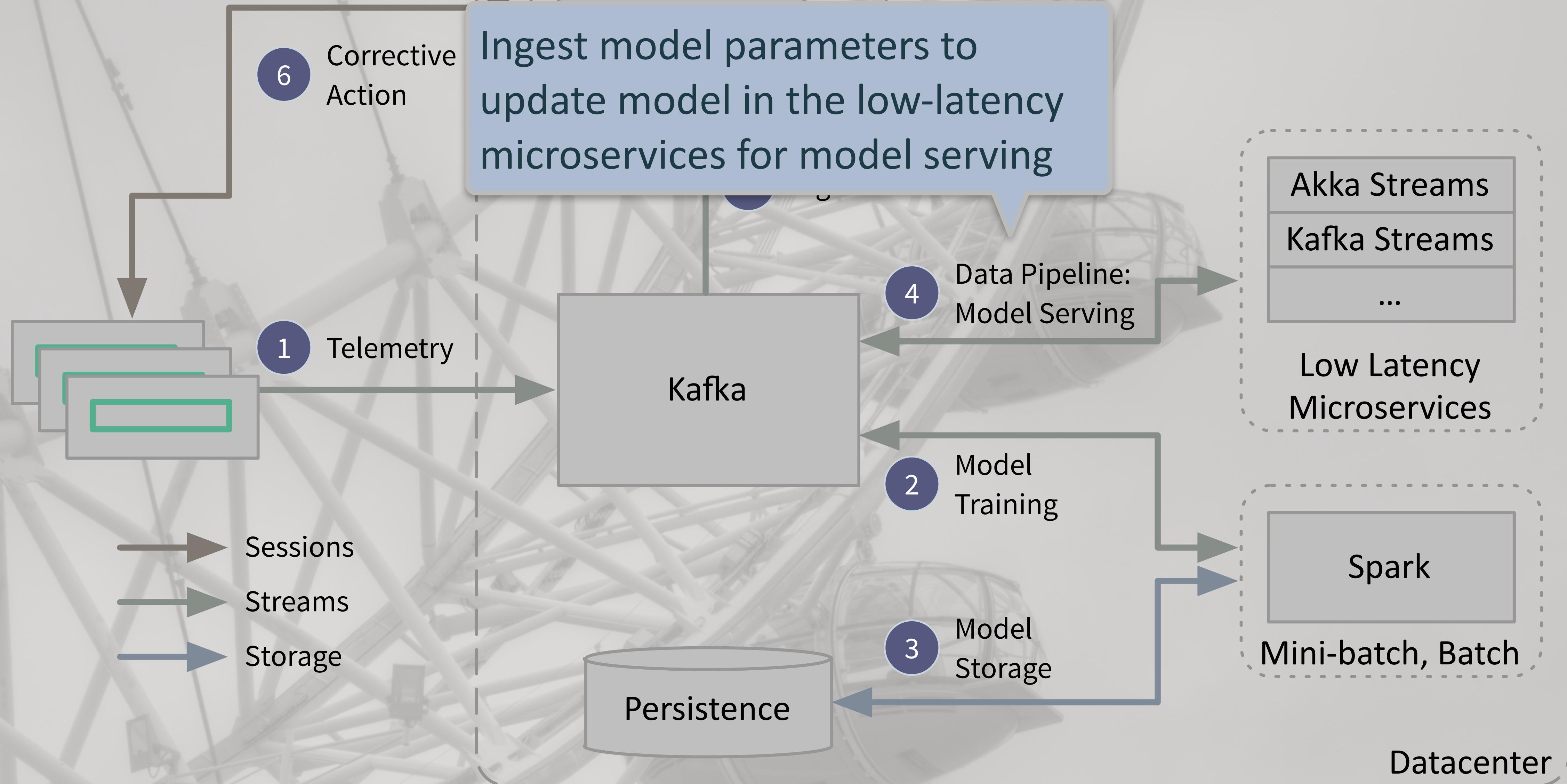
Example Architecture



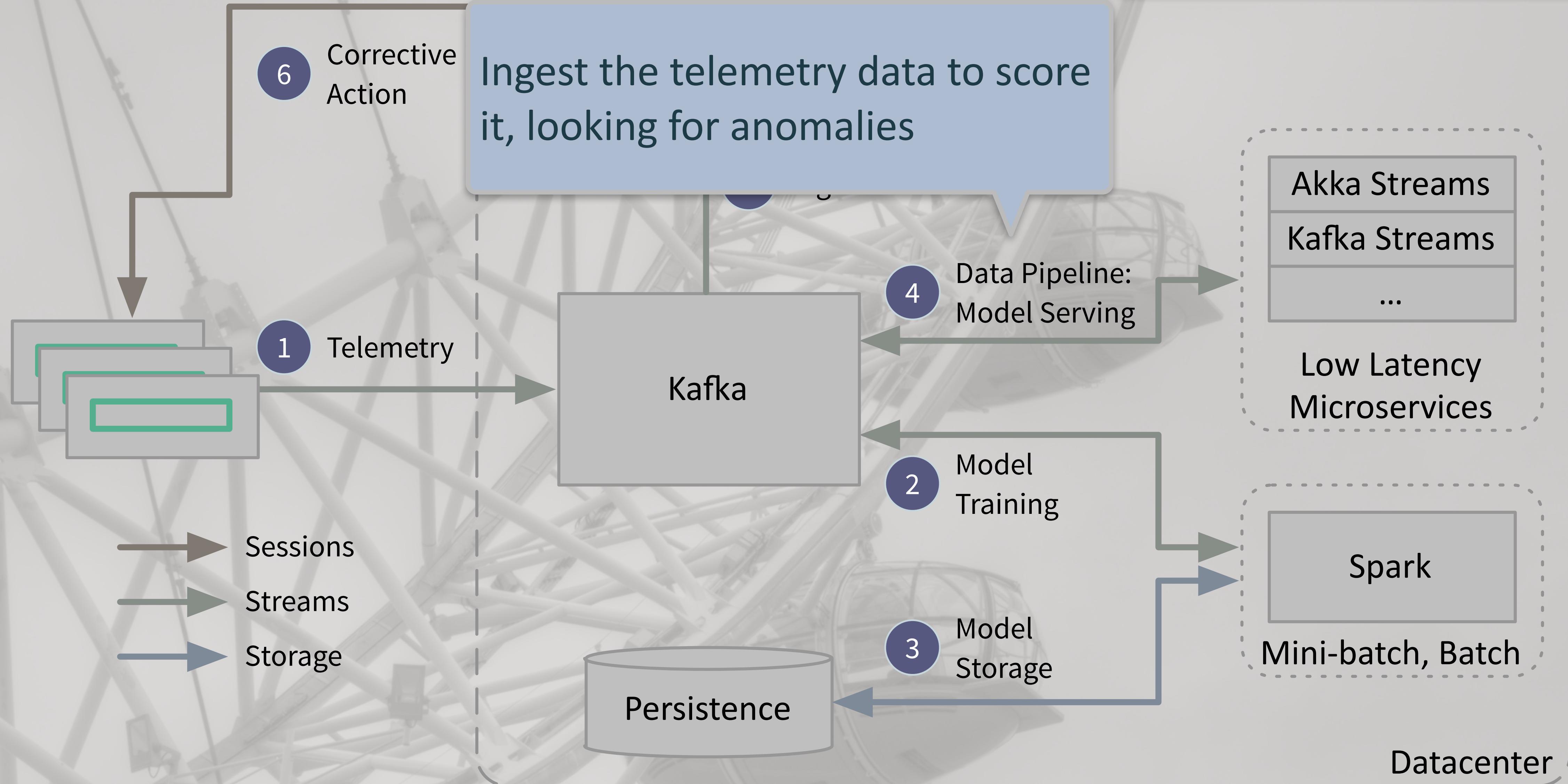
Example Architecture



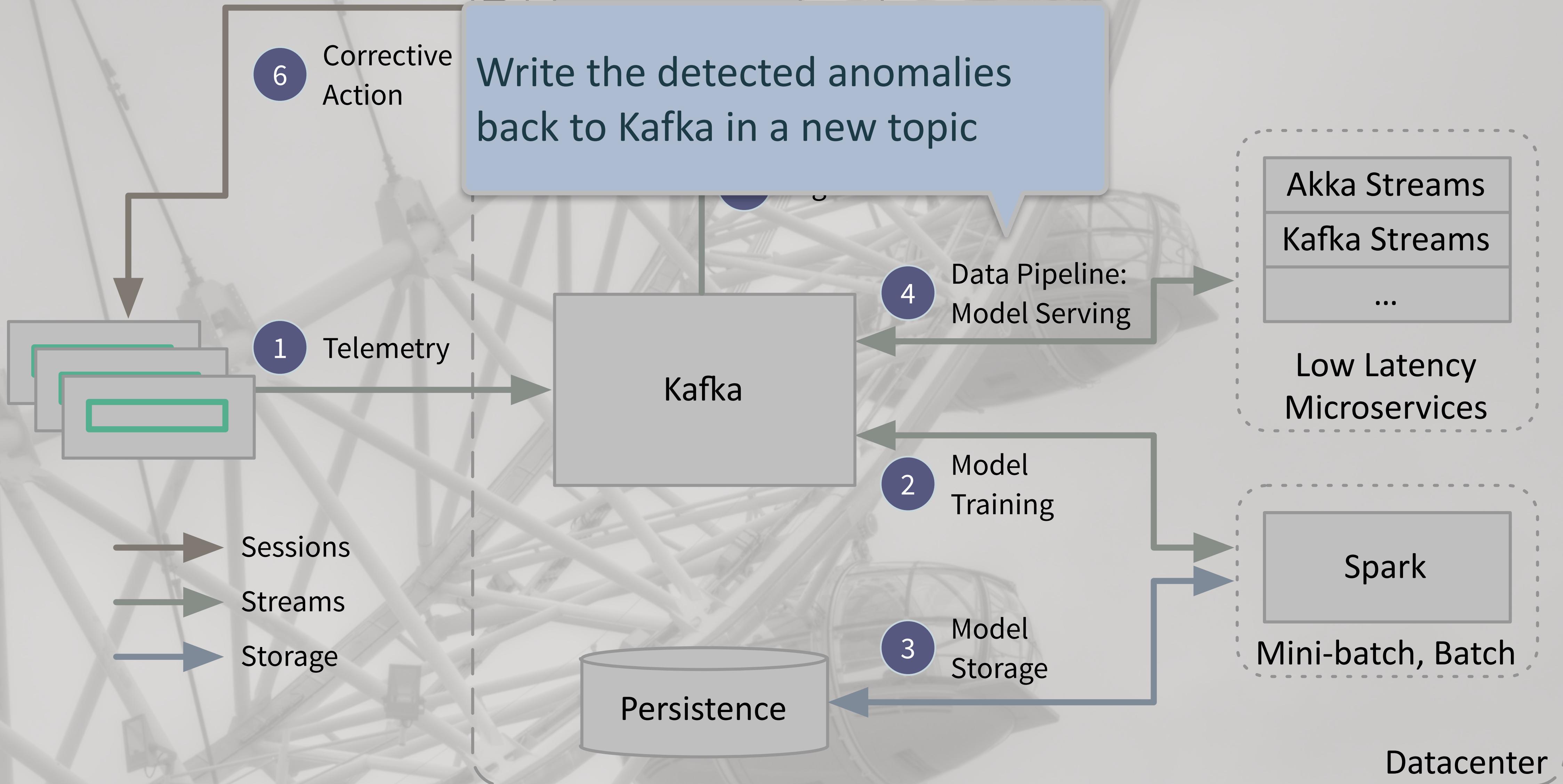
Example Architecture



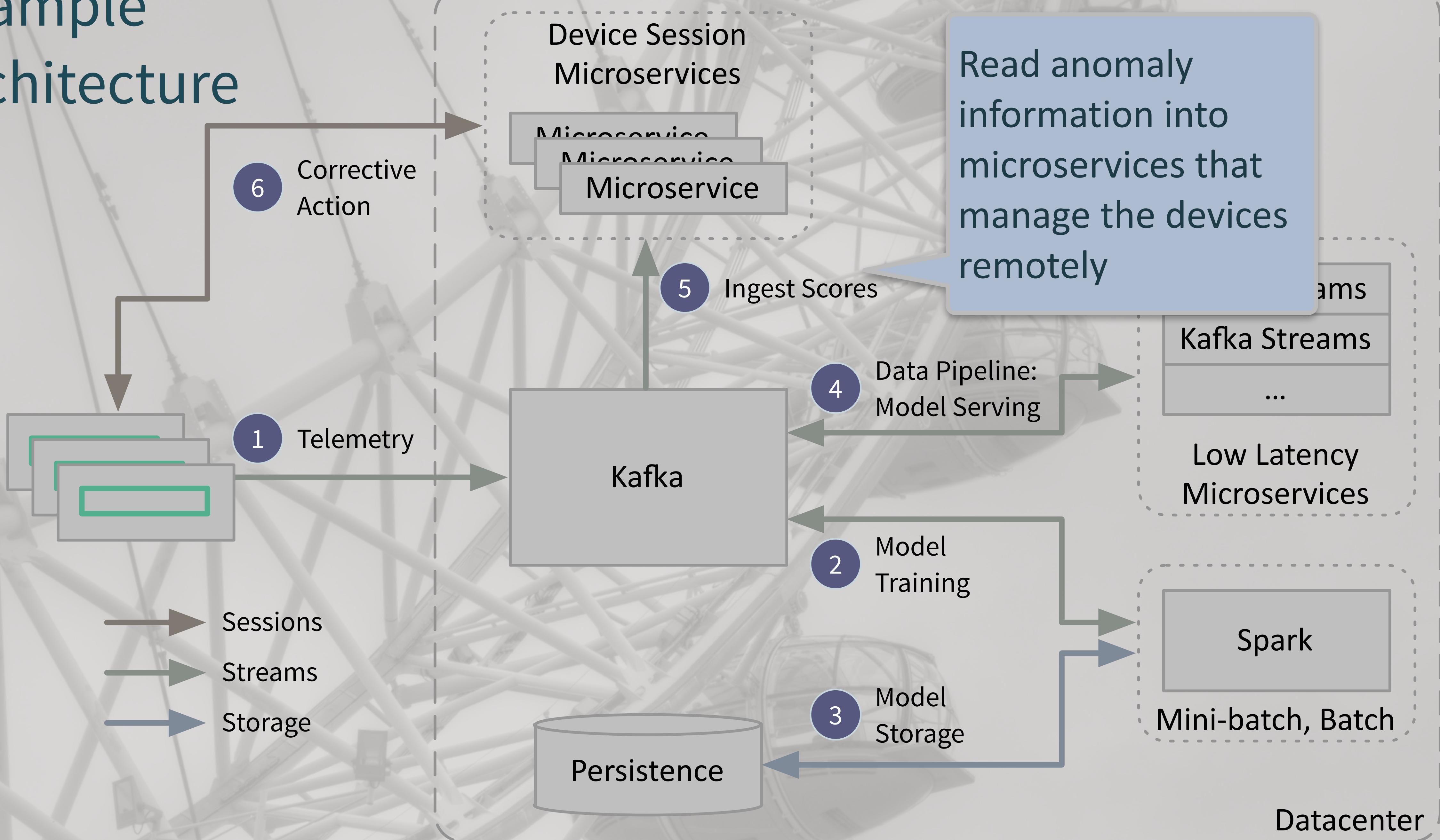
Example Architecture



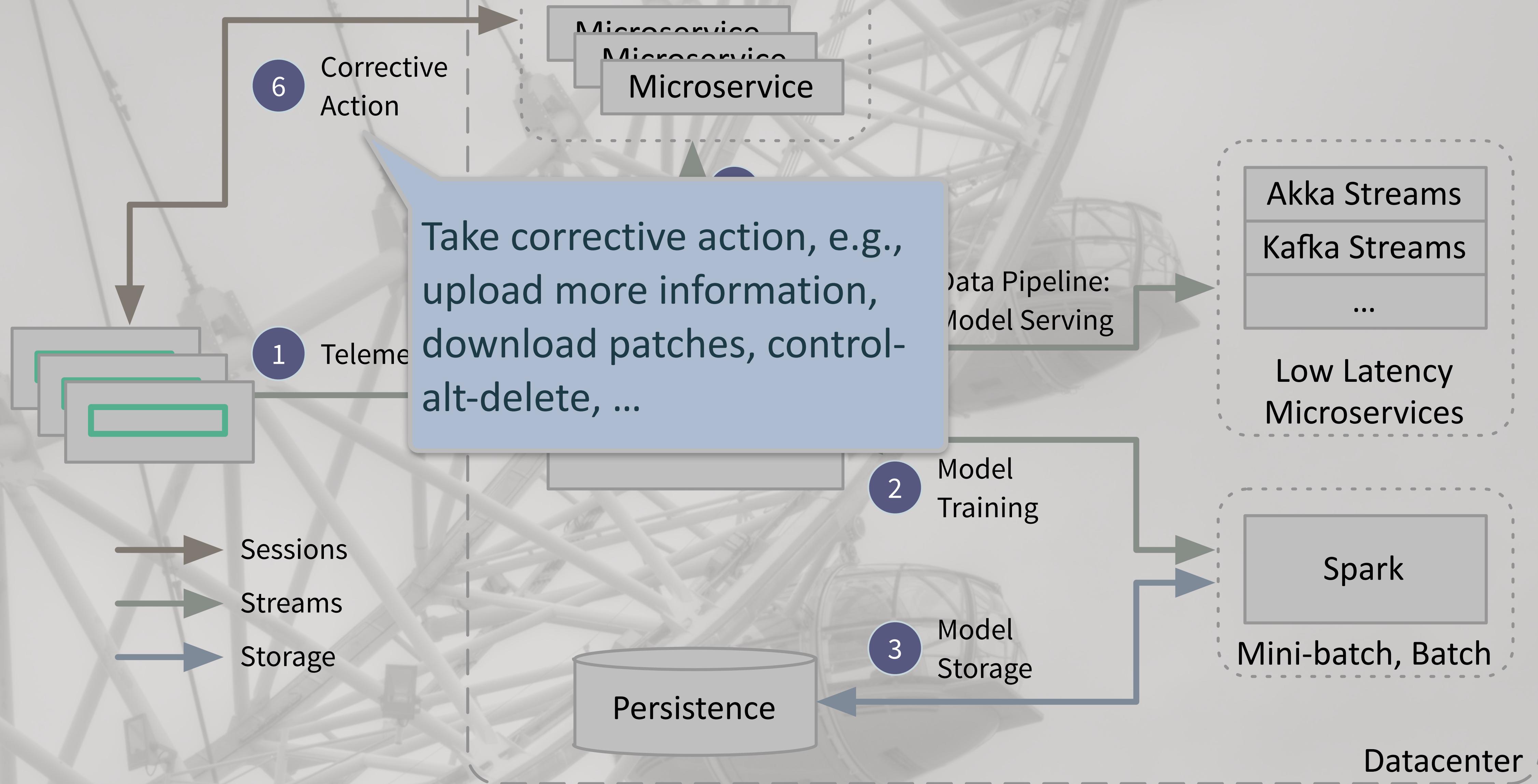
Example Architecture



Example Architecture

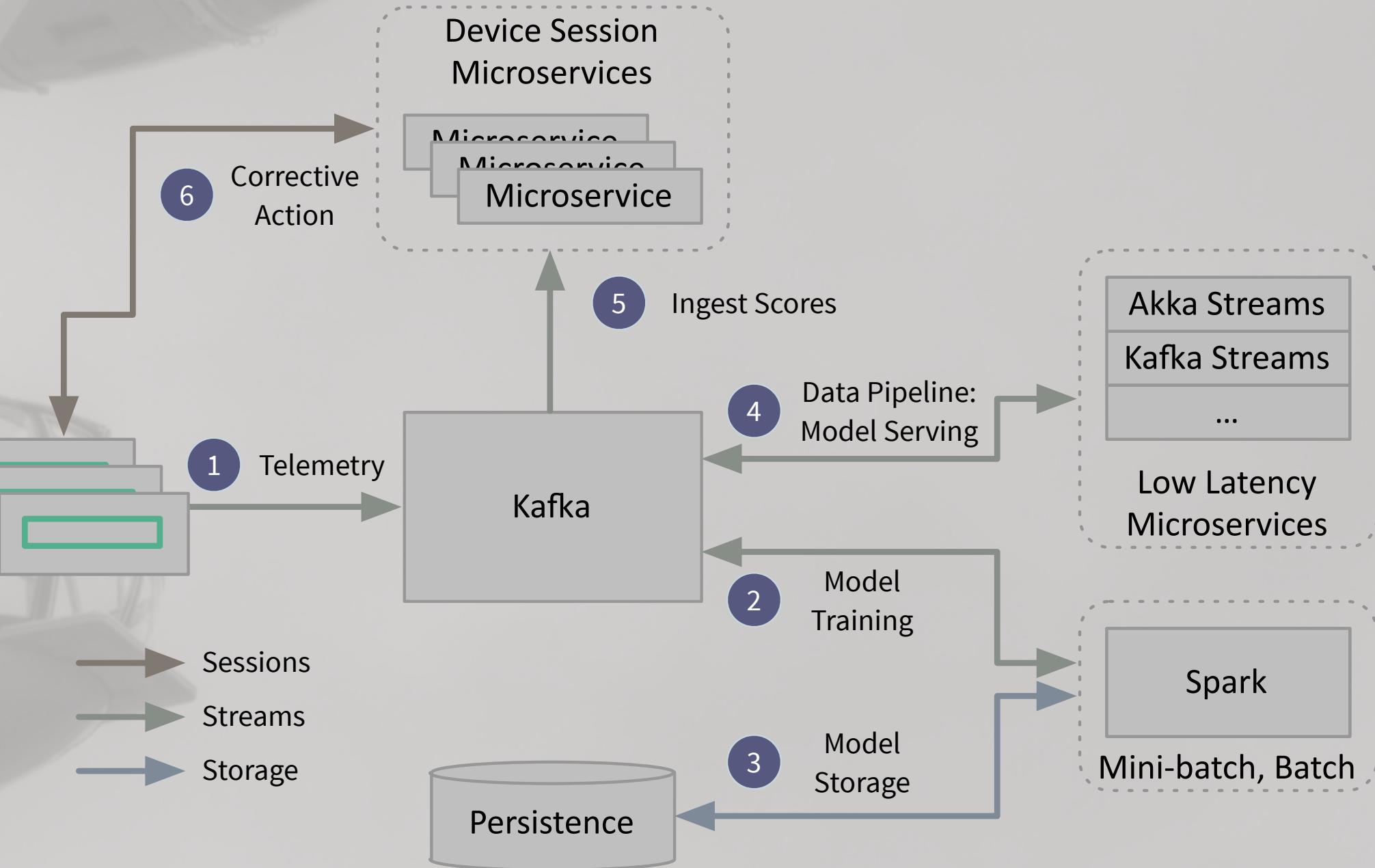


Example Architecture



Challenges

- Network overhead for telemetry ingestion too high?
- Model serving latency too long?
- Idea: Serve model on the device!

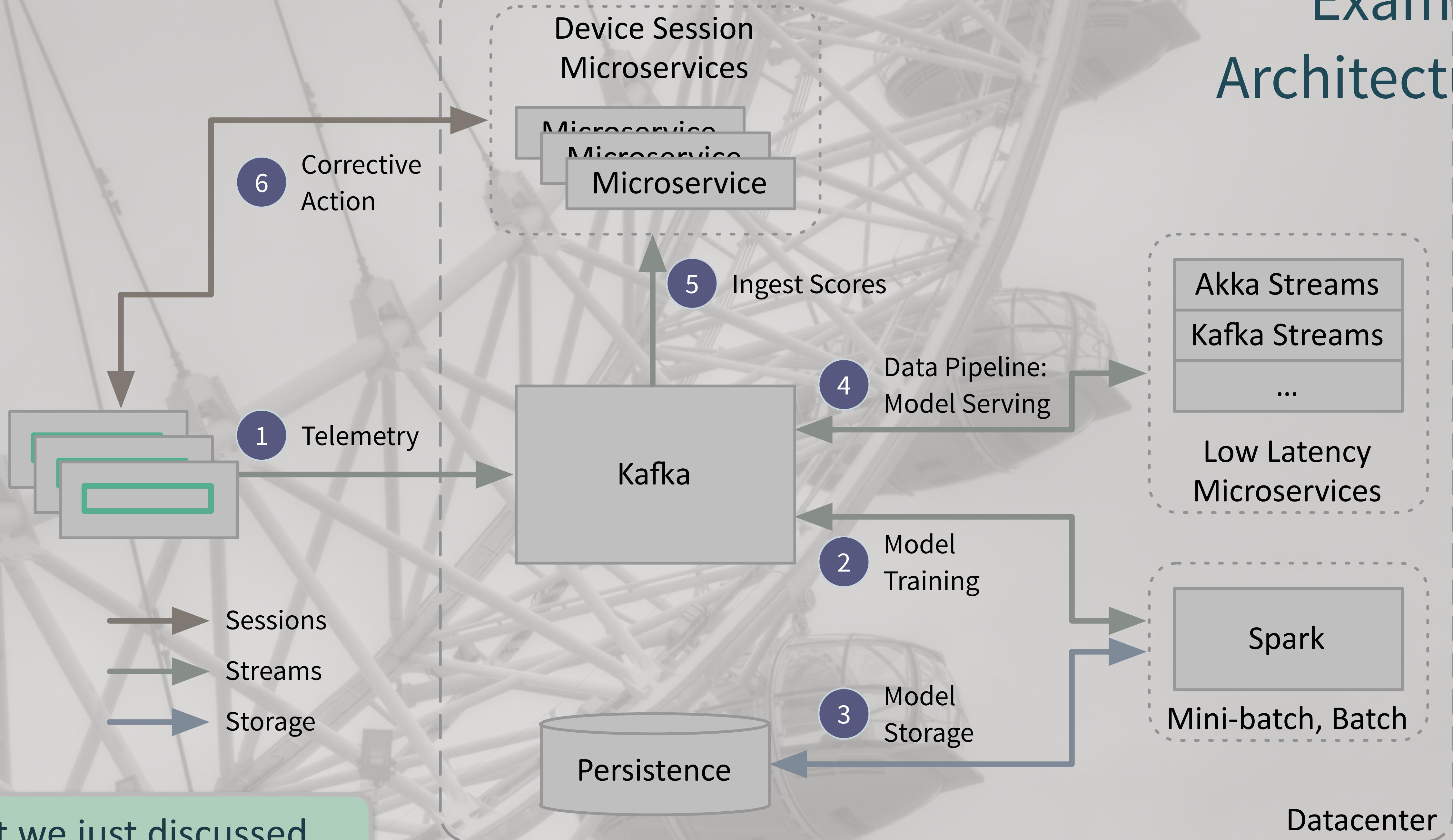




Internet of Things

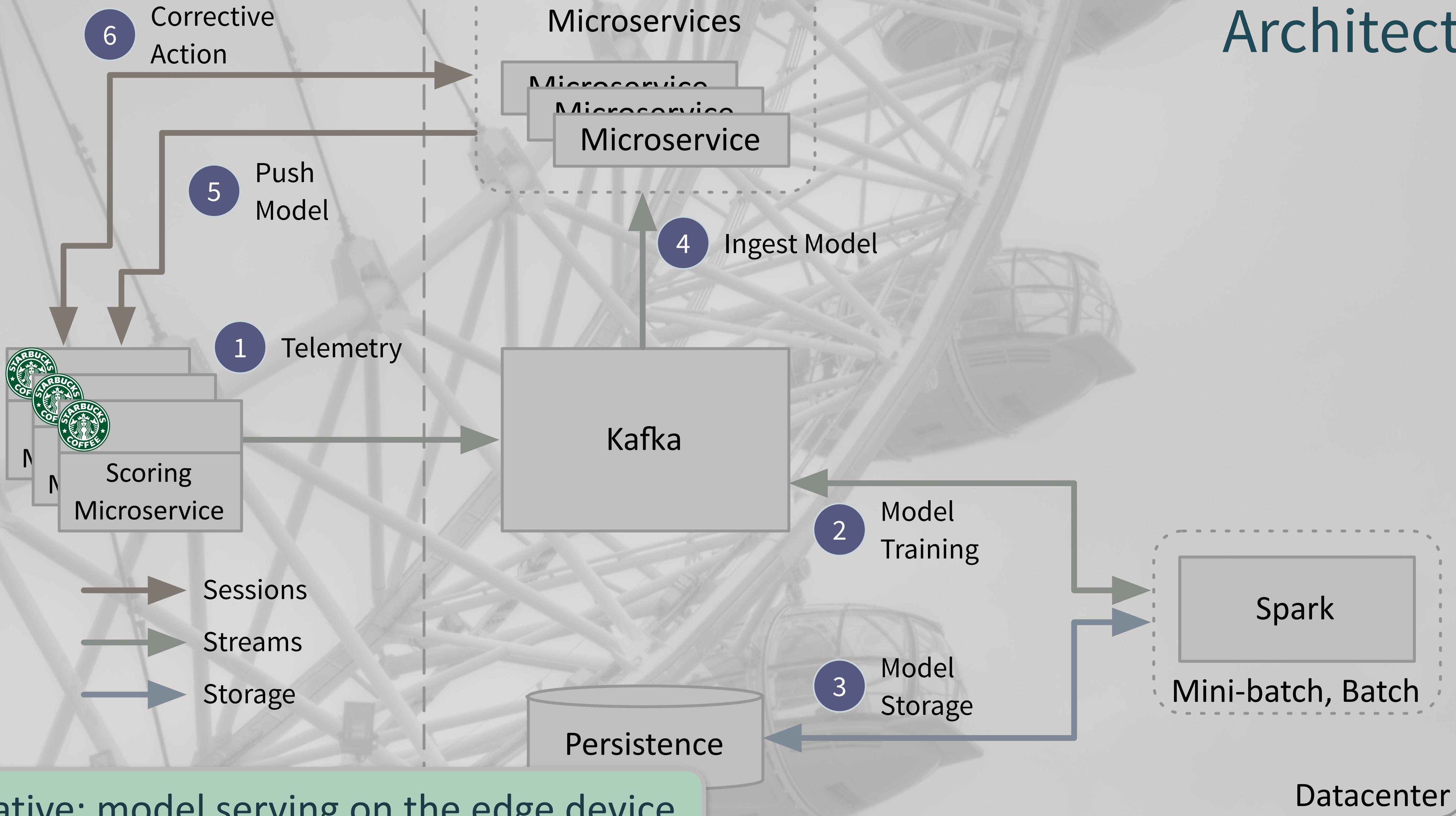
- Real-time consumer and industrial device and supply chain management at scale

Example Architecture

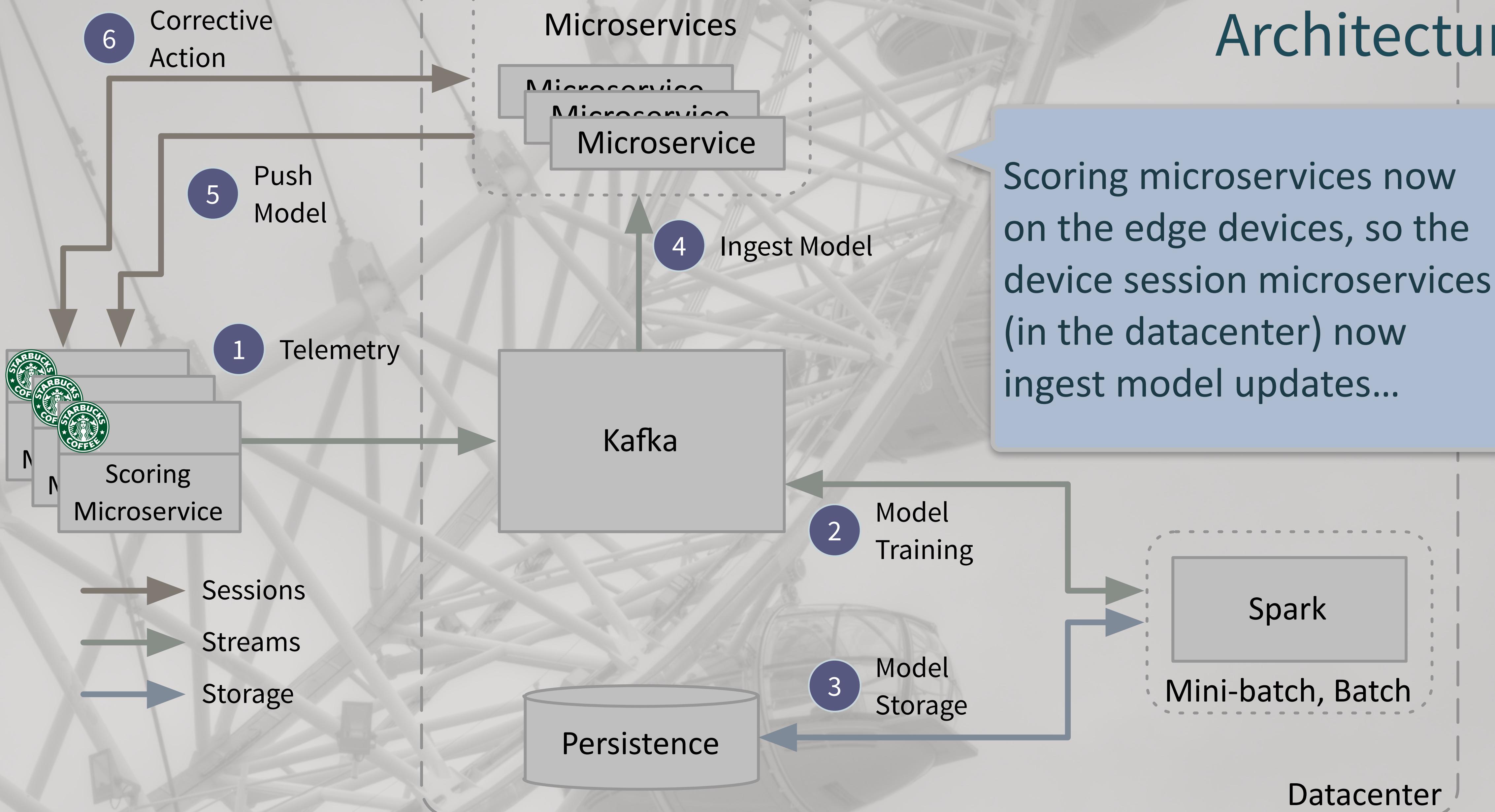


What we just discussed...

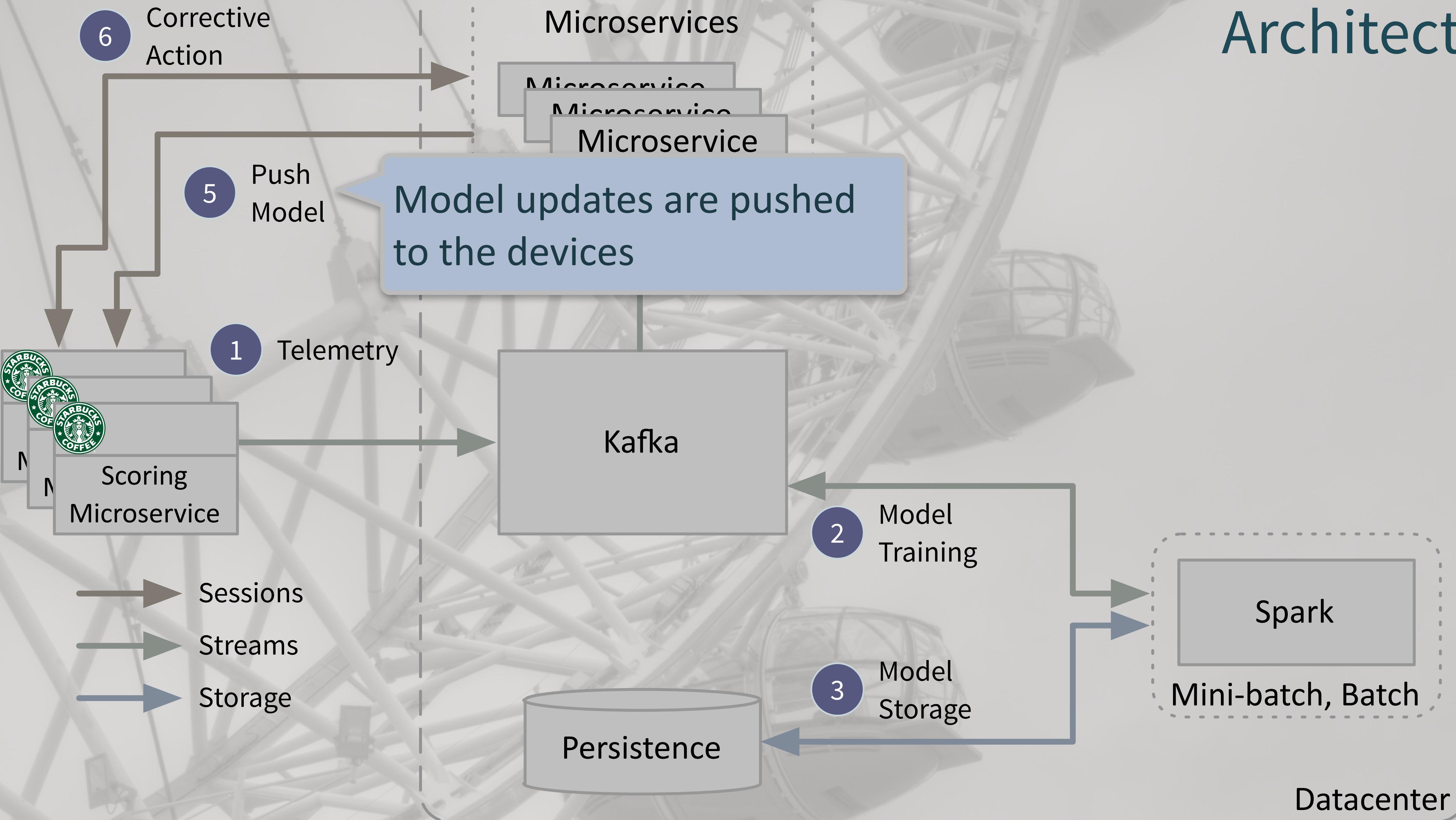
Edge Scoring Example Architecture



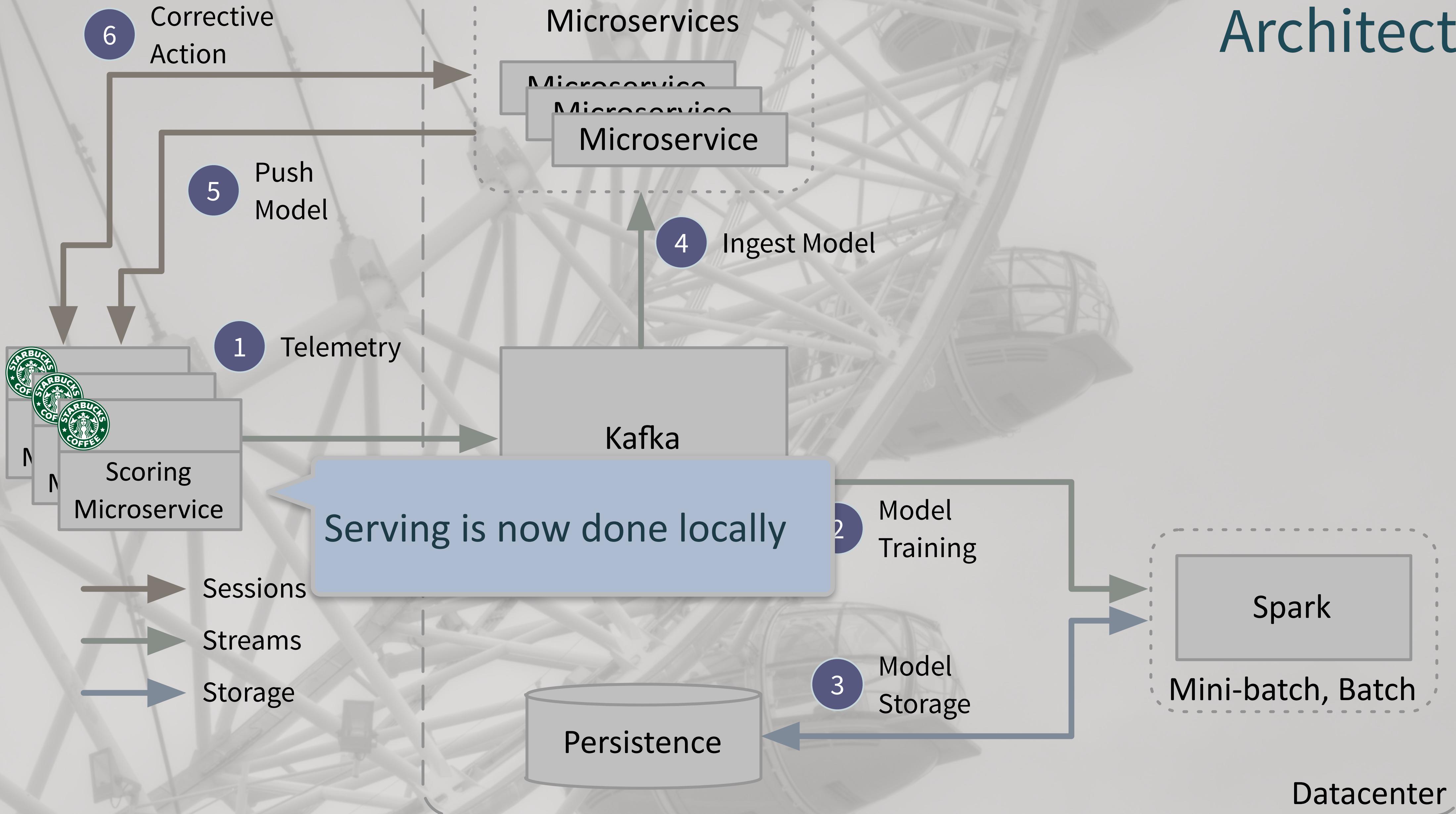
Edge Scoring Example Architecture



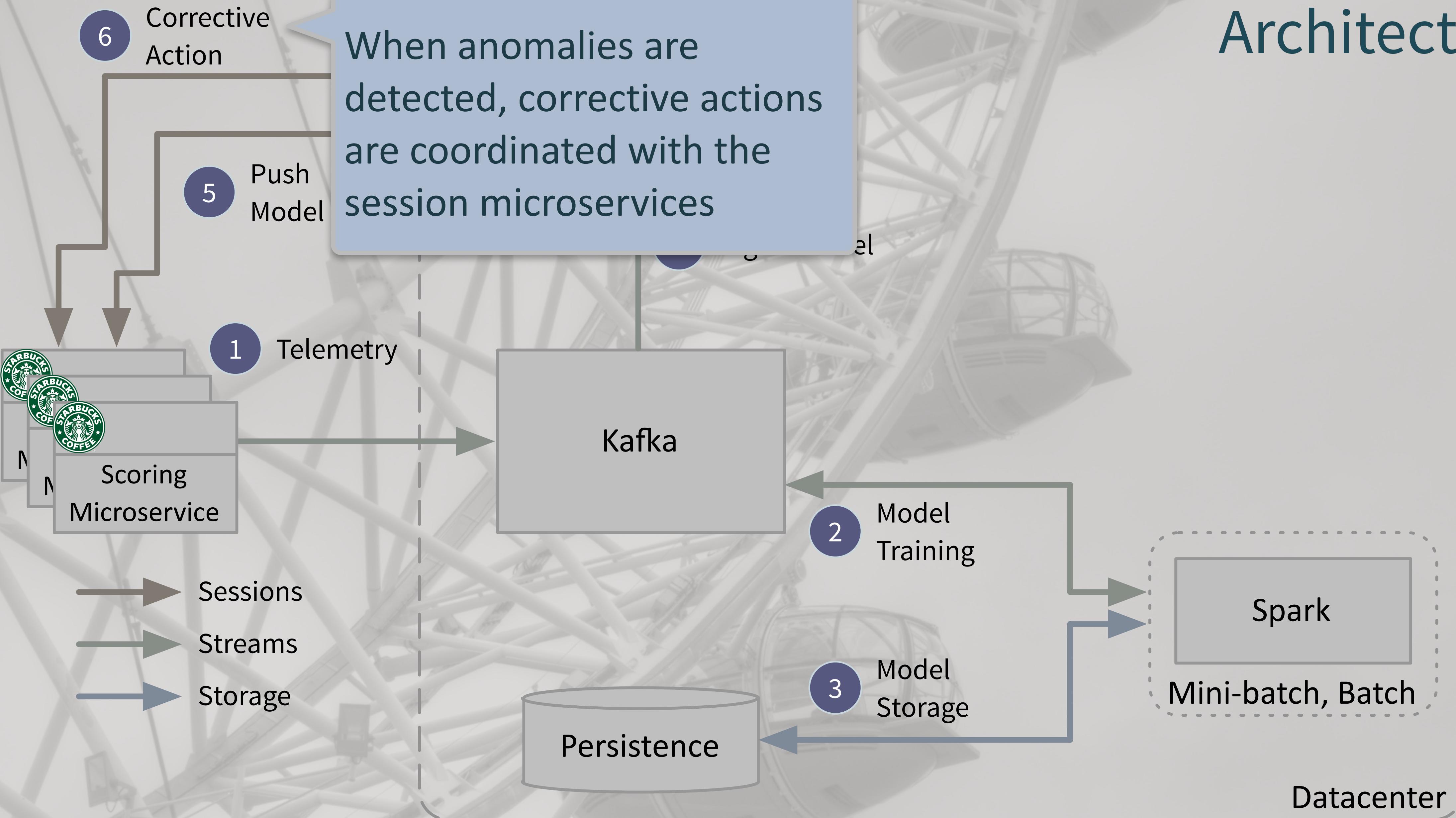
Edge Scoring Example Architecture



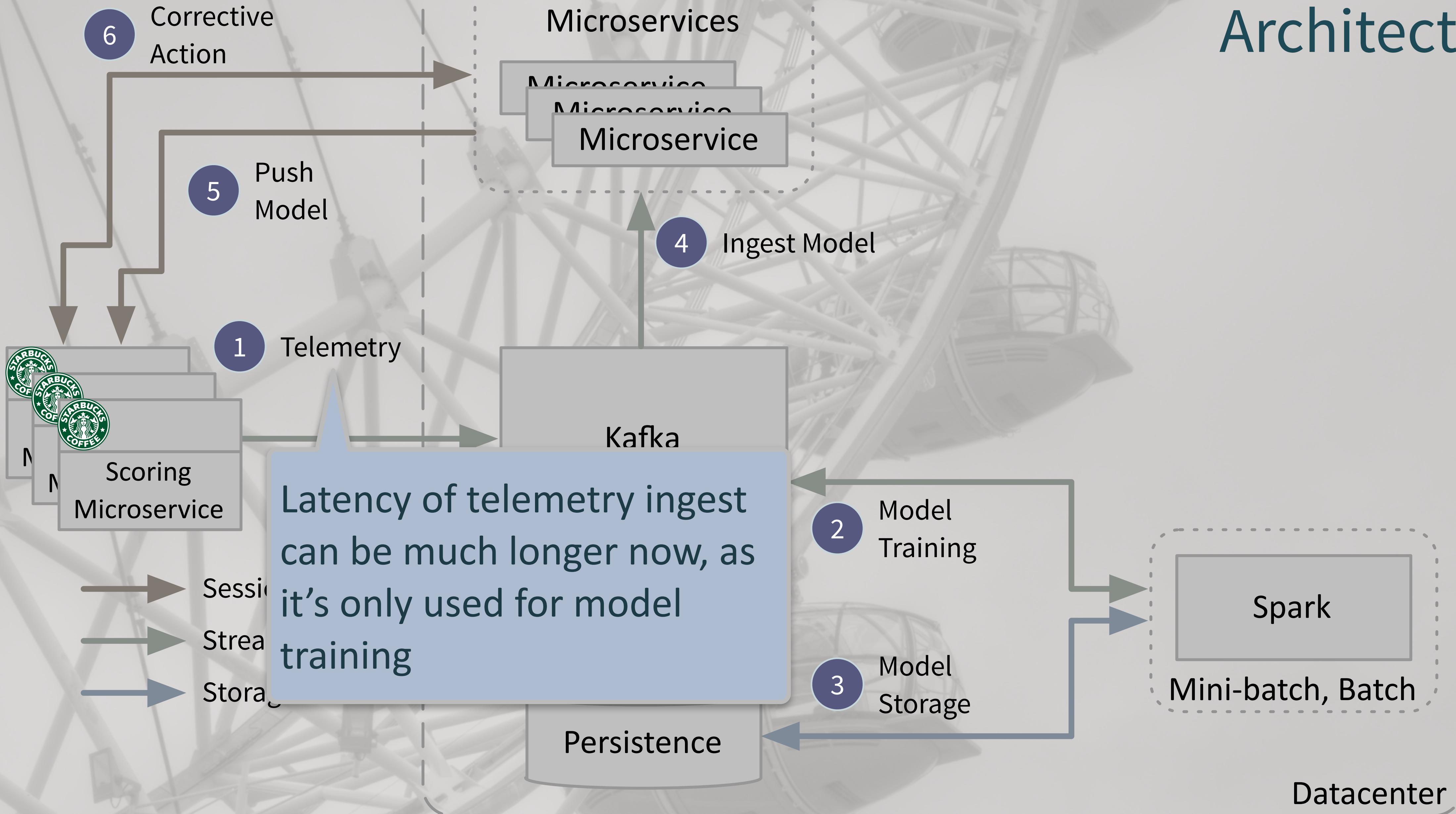
Edge Scoring Example Architecture



Edge Scoring Example Architecture



Edge Scoring Example Architecture



Fas

Batch -> streaming for competitive advantage

Cases

Predictive Analytics

Apply ML models to large volumes of device data to pre-empt failures / outages

 Hewlett Packard Enterprise

IoT

Real-time consumer and industrial Device and Supply Chain management at scale



Real-time Personalization

Real-time marketing based on behavior, location, inventory levels, product promotions, etc.



Real-time Financial Processes

Drive better business outcomes through real-time risk, fraud detection, compliance, audit, governance, etc.



Legacy Modernization

Accelerate decision making processes and optimize infrastructure costs by moving from batch to streaming



More at: <https://www.lightbend.com/customers>

Real-time Personalization



- Model users to provide real-time marketing based on behavior, location, inventory levels, product promotions, etc.

Real-time Financial



- Drive better business outcomes through real-time risk, fraud detection, compliance, audit, governance, etc.

Legacy Modernization



- Accelerate decision making processes and optimize infrastructure costs by moving from batch to streaming
- Hadoop replacement



Technology Choices

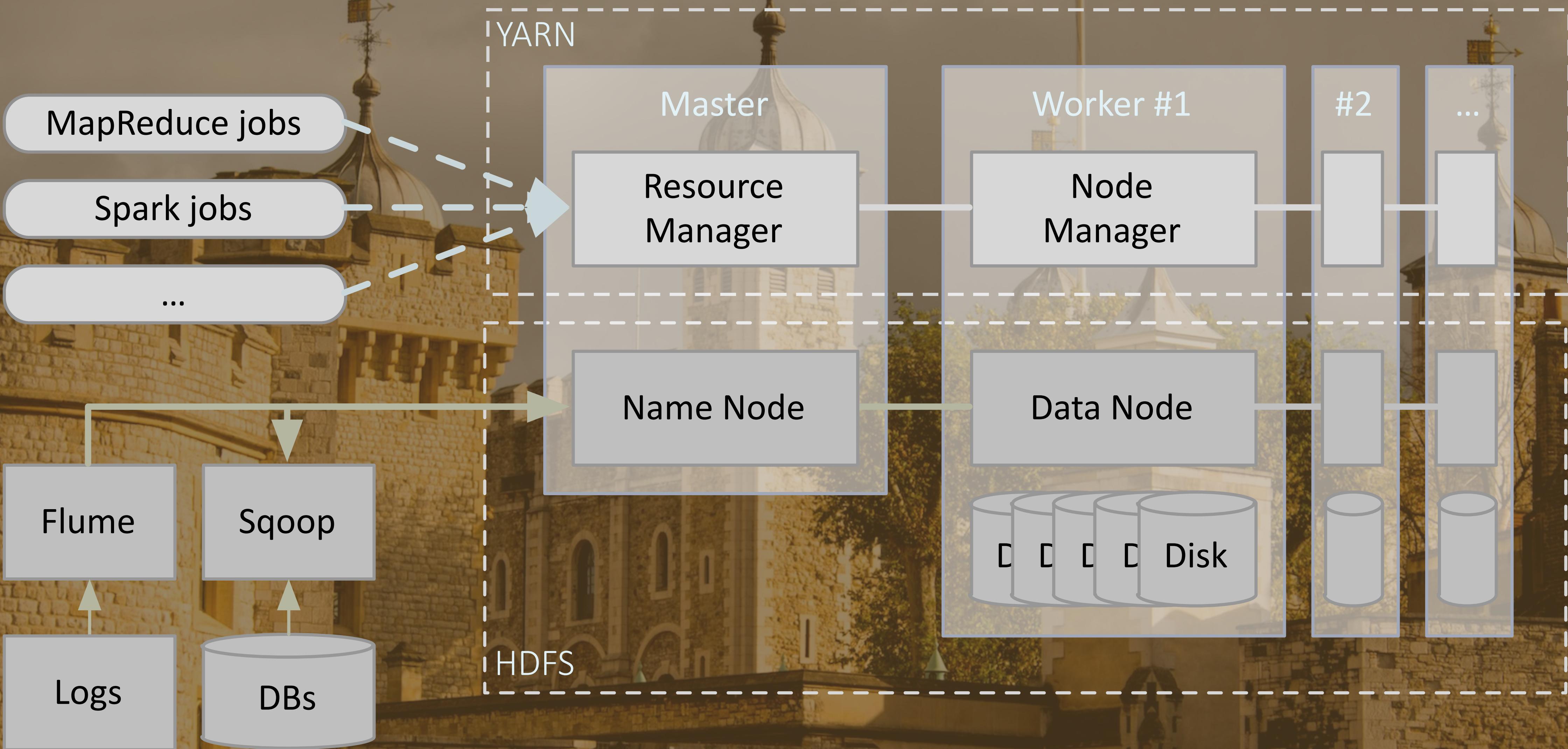
Technology Choices

- More than “tweaking” Hadoop...
- New architectures that merge data processing with microservices

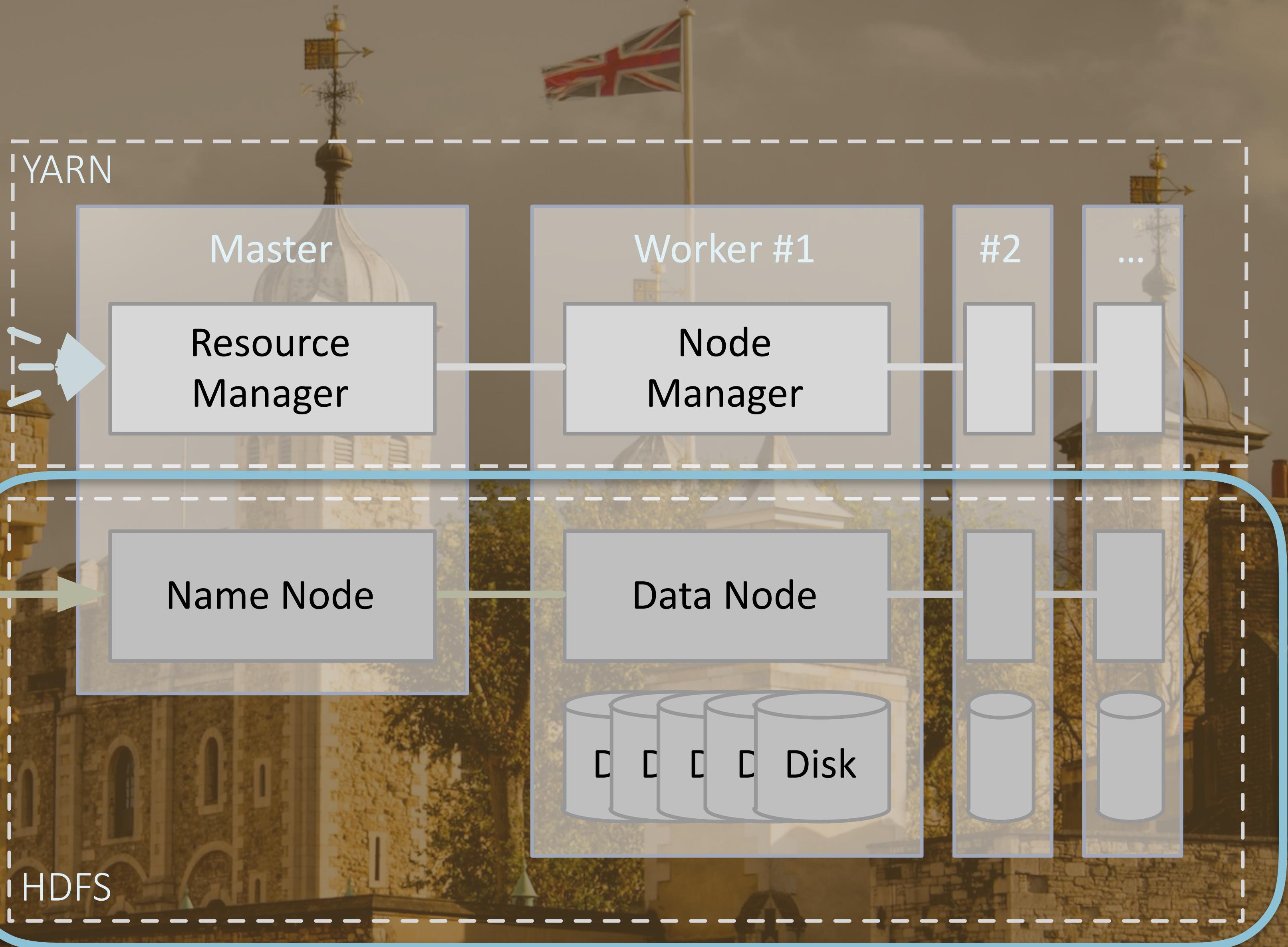
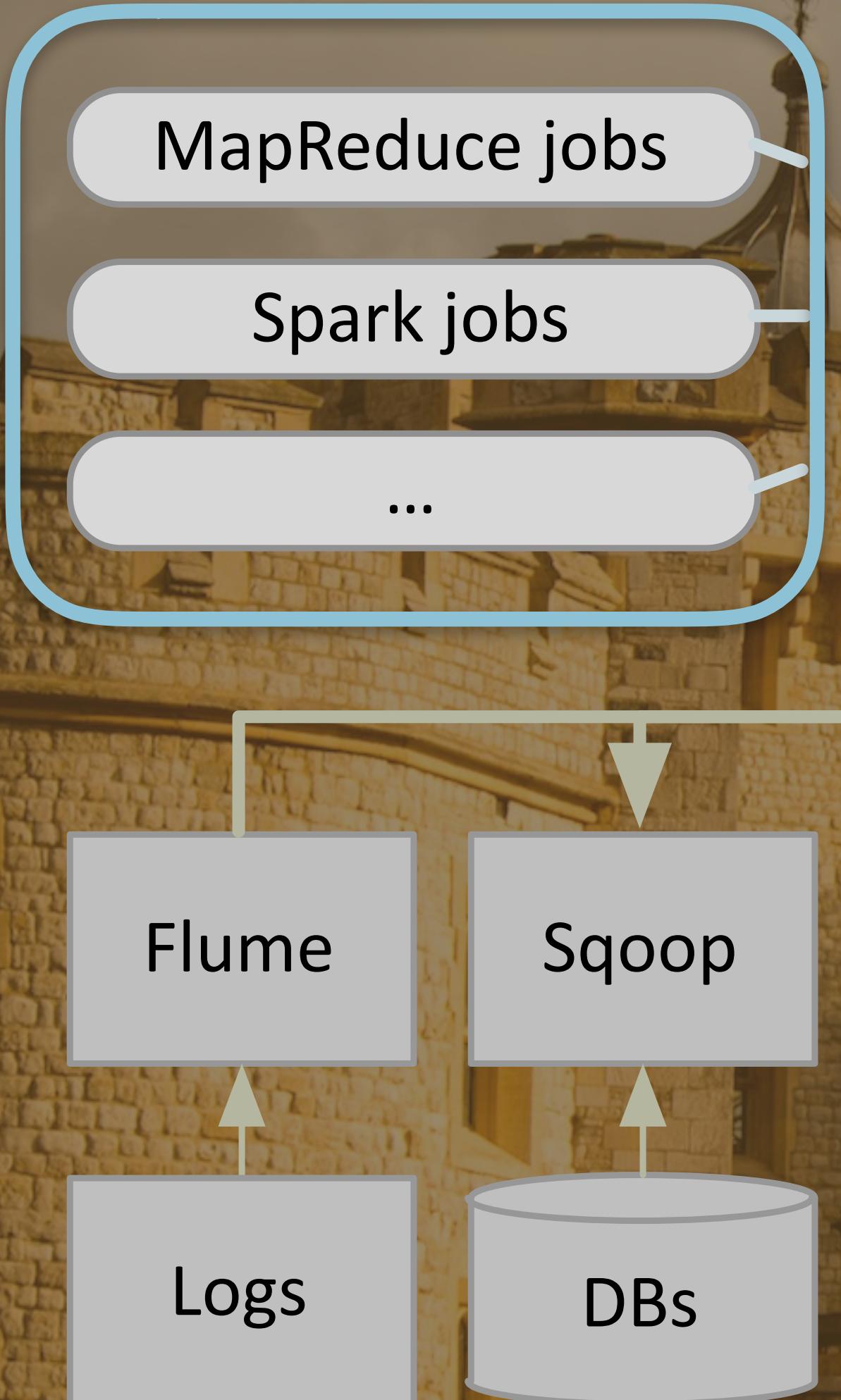
Recall Hadoop...



- 
- A photograph of a historic stone building, likely a university or institutional complex, featuring multiple towers and domes. A Union Jack flag flies from a pole on top of one of the towers. The building is made of light-colored stone and has several arched windows and doorways. The sky is overcast.
- Data warehouse replacement
 - Historical analysis
 - Interactive exploration
 - Offline training of machine learning models
 - ...

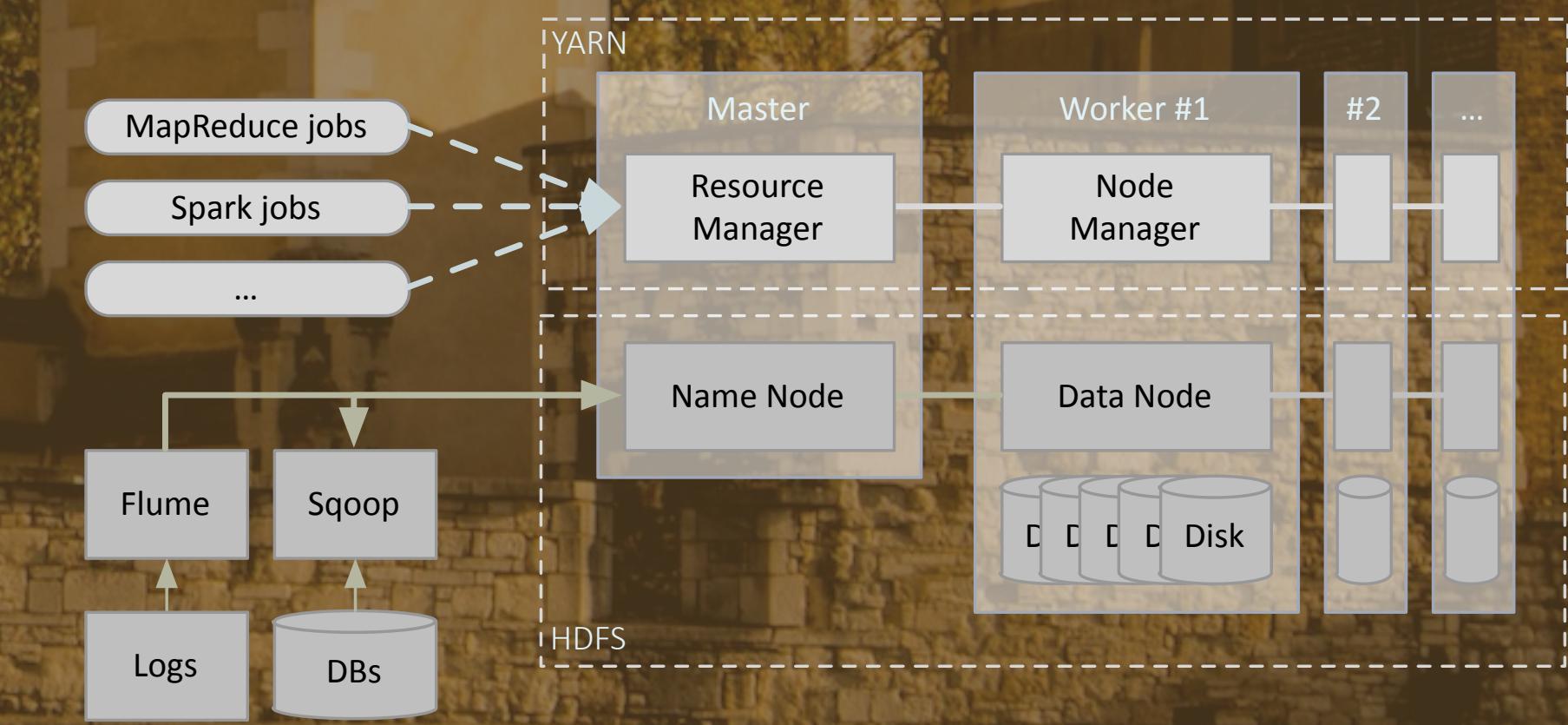


Compute



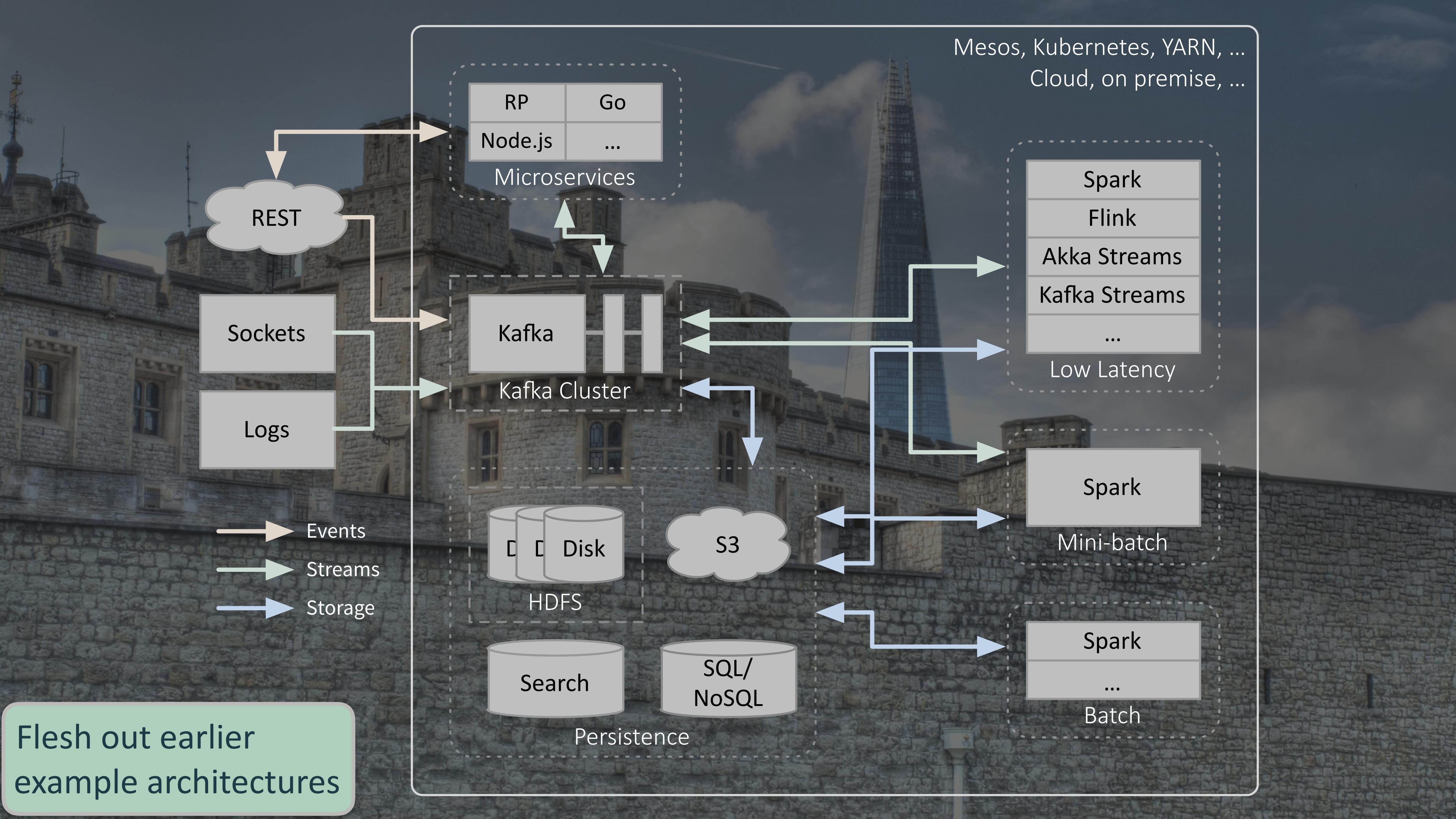
Storage

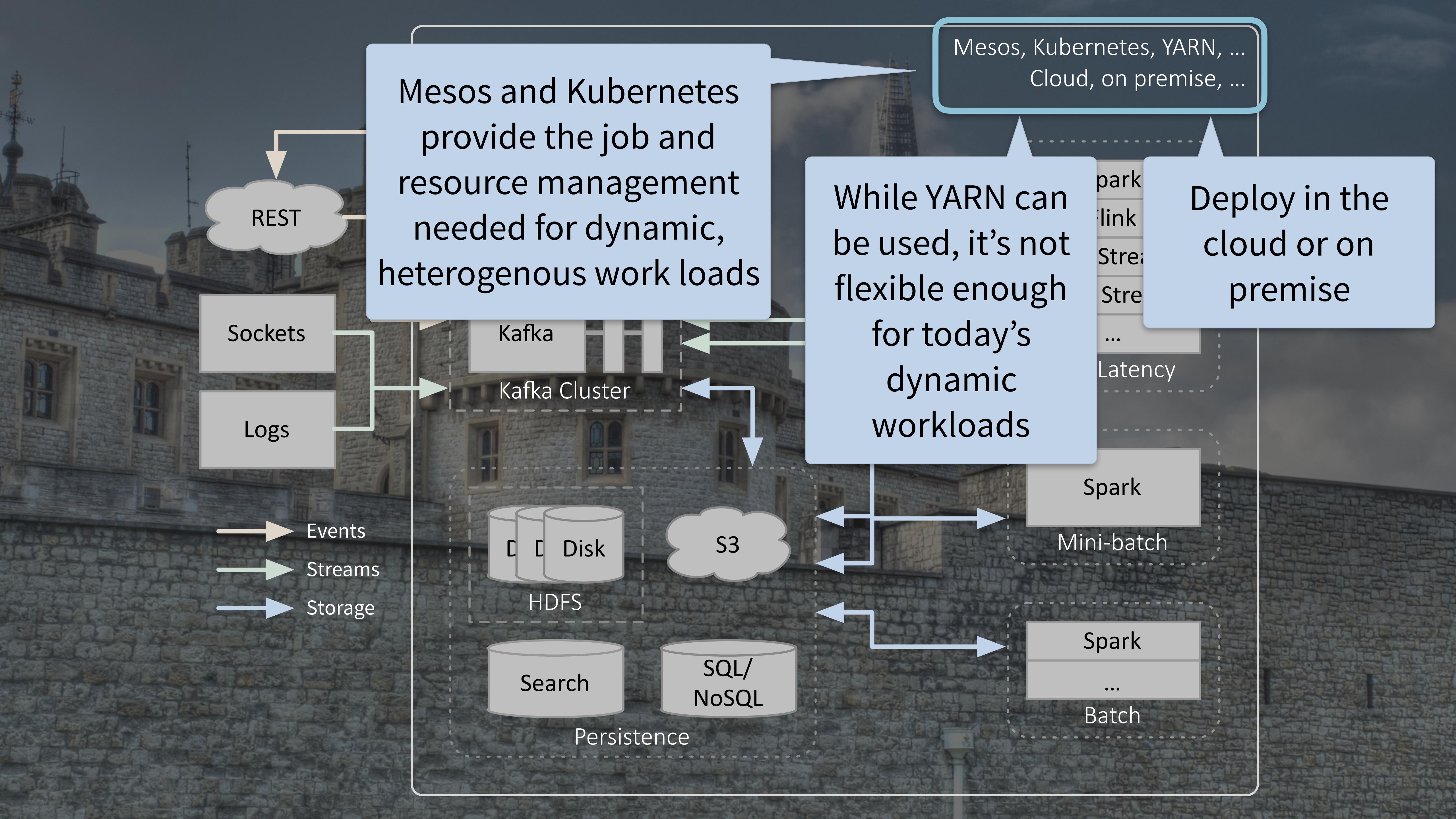
- Hadoop is ideal for batch and interactive apps
- ... but also constrained by that model





New Fast Data Architecture



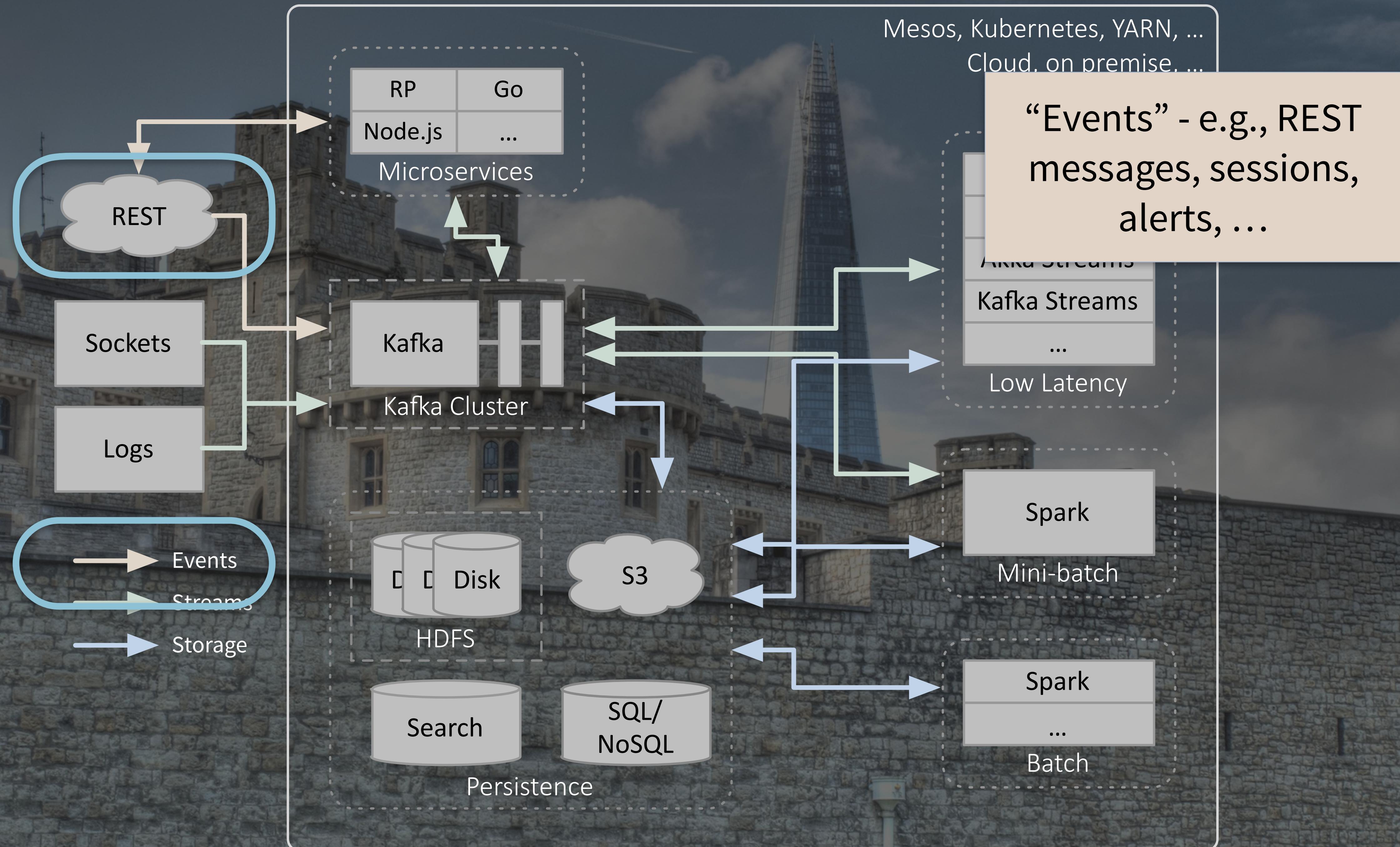


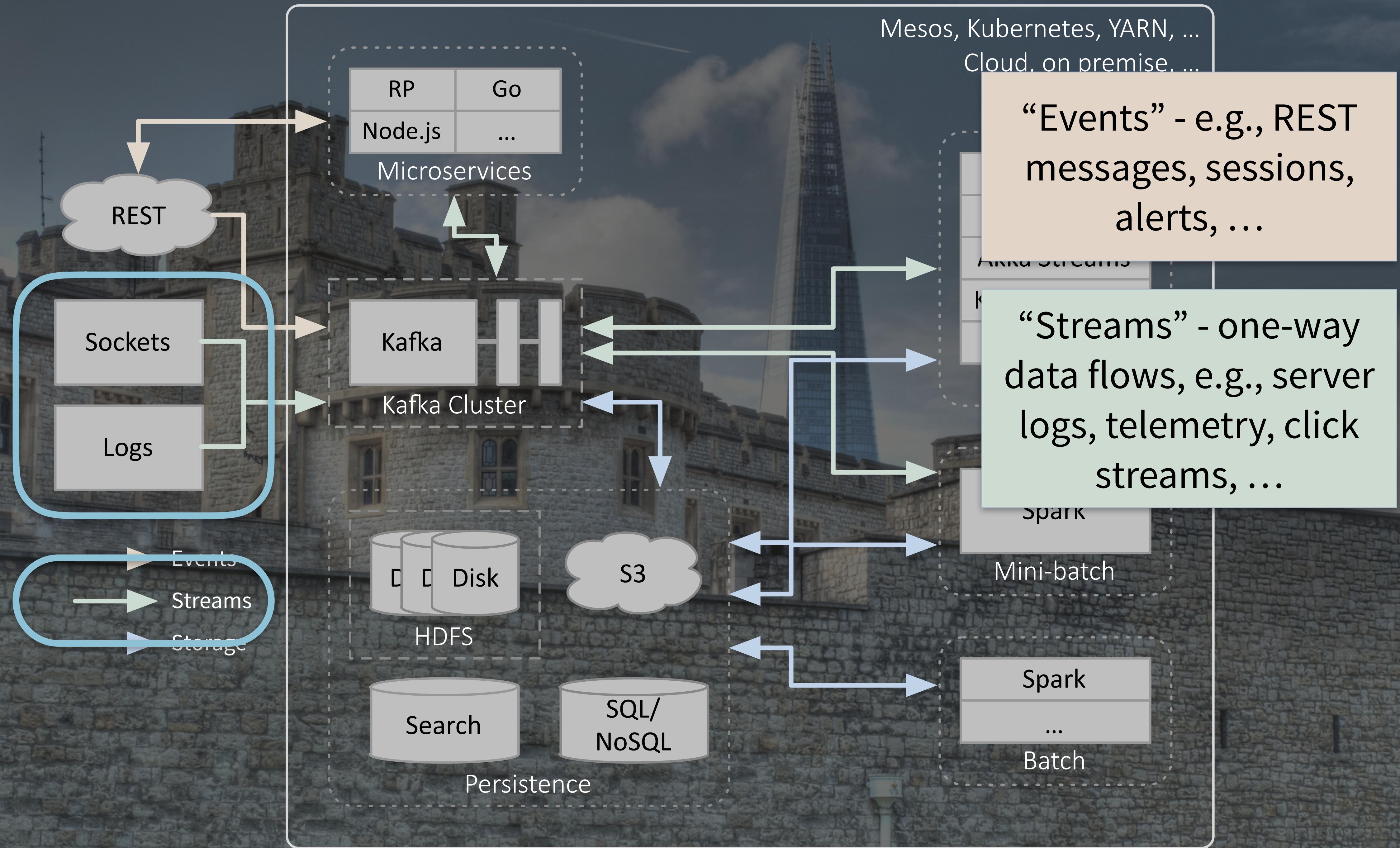
Mesos and Kubernetes provide the job and resource management needed for dynamic, heterogenous work loads

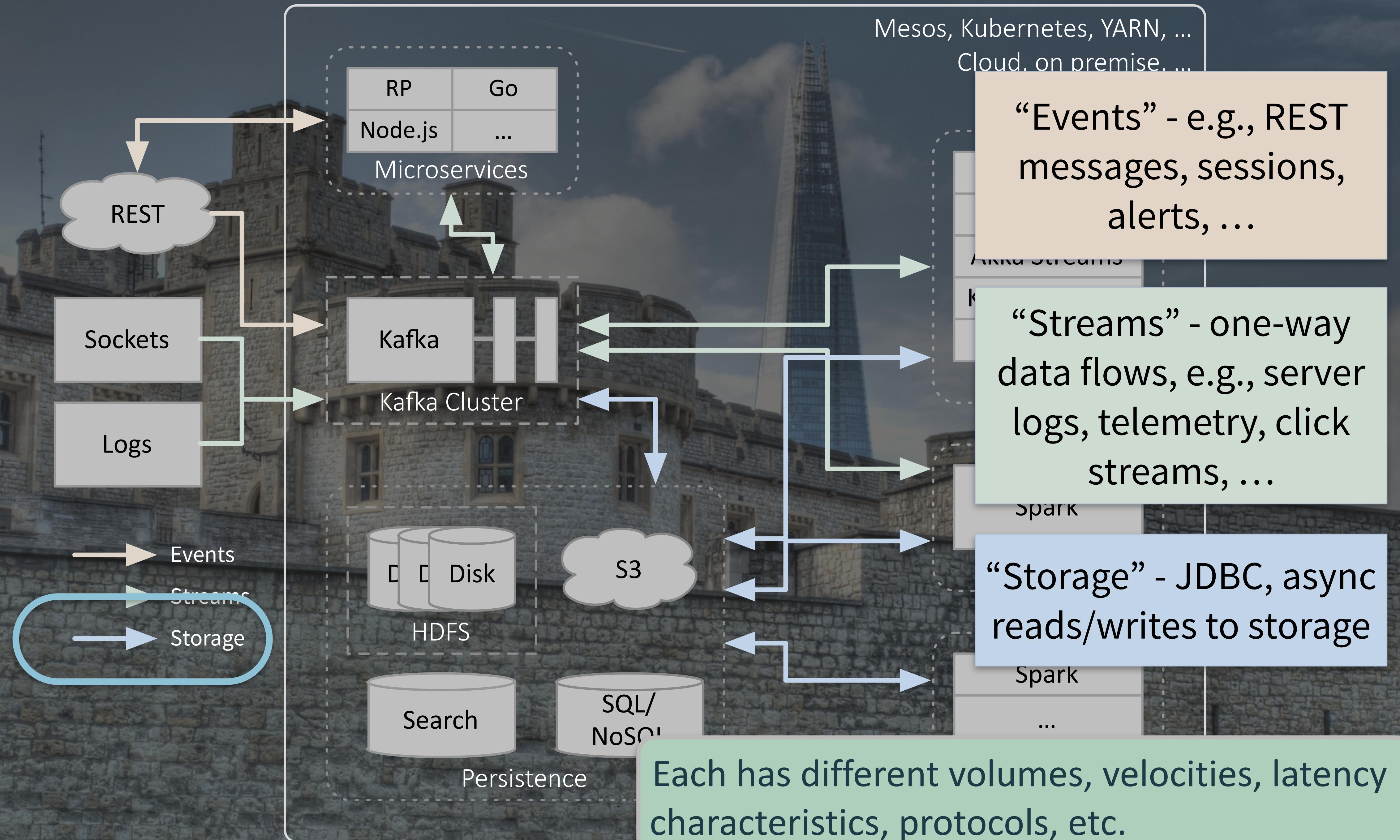
While YARN can be used, it's not flexible enough for today's dynamic workloads

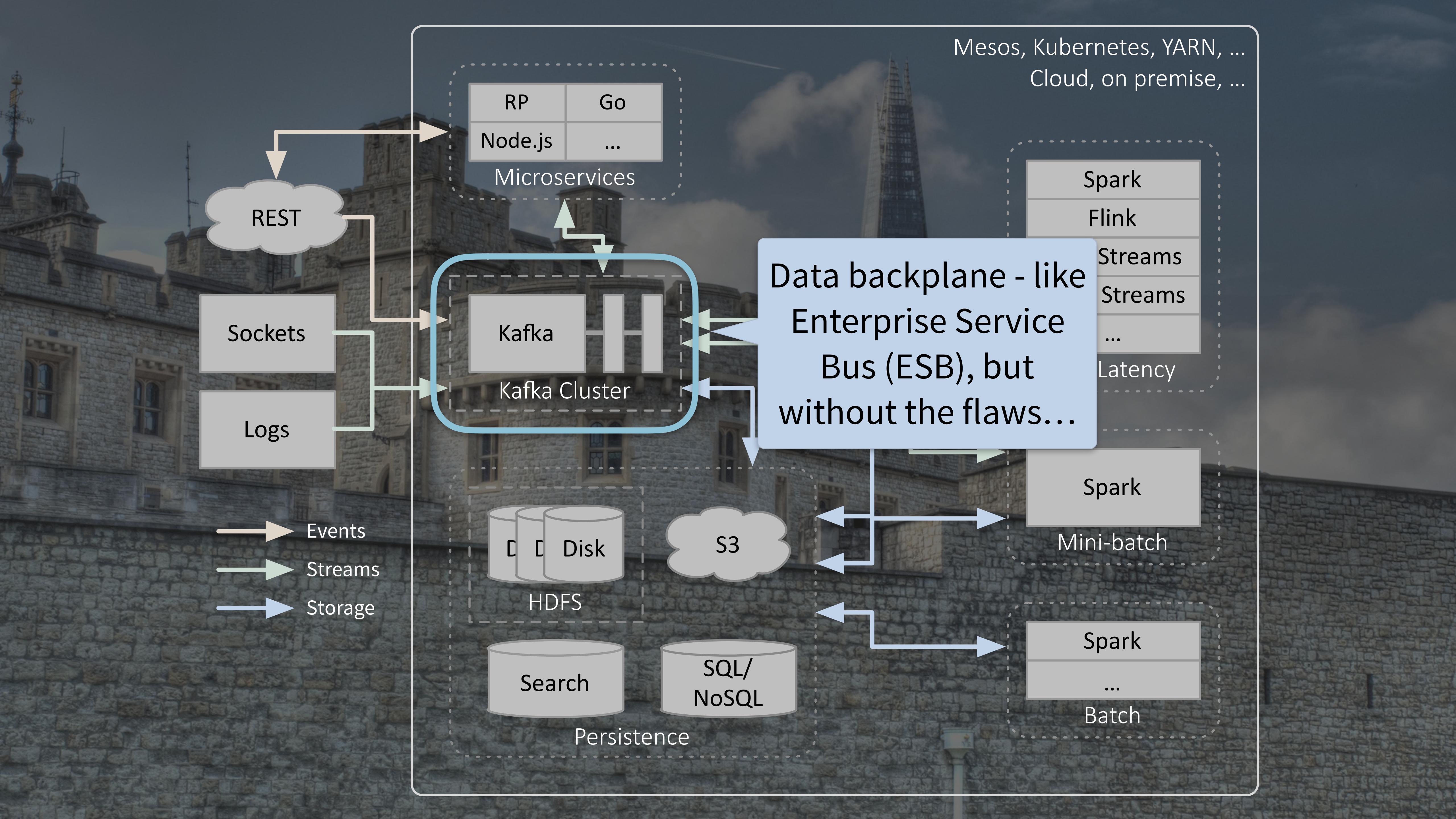
Mesos, Kubernetes, YARN, ...
Cloud, on premise, ...

Deploy in the cloud or on premise



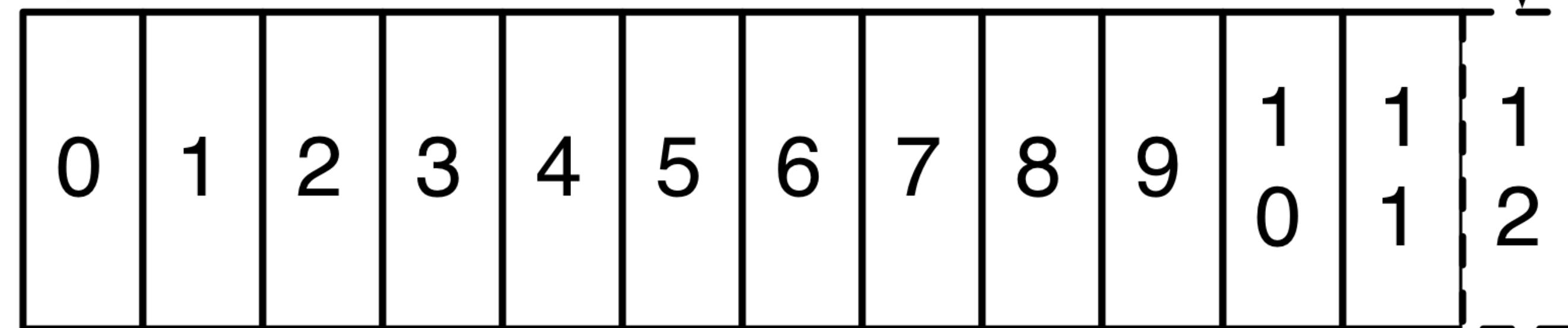






Why Kafka?

Organized into topics

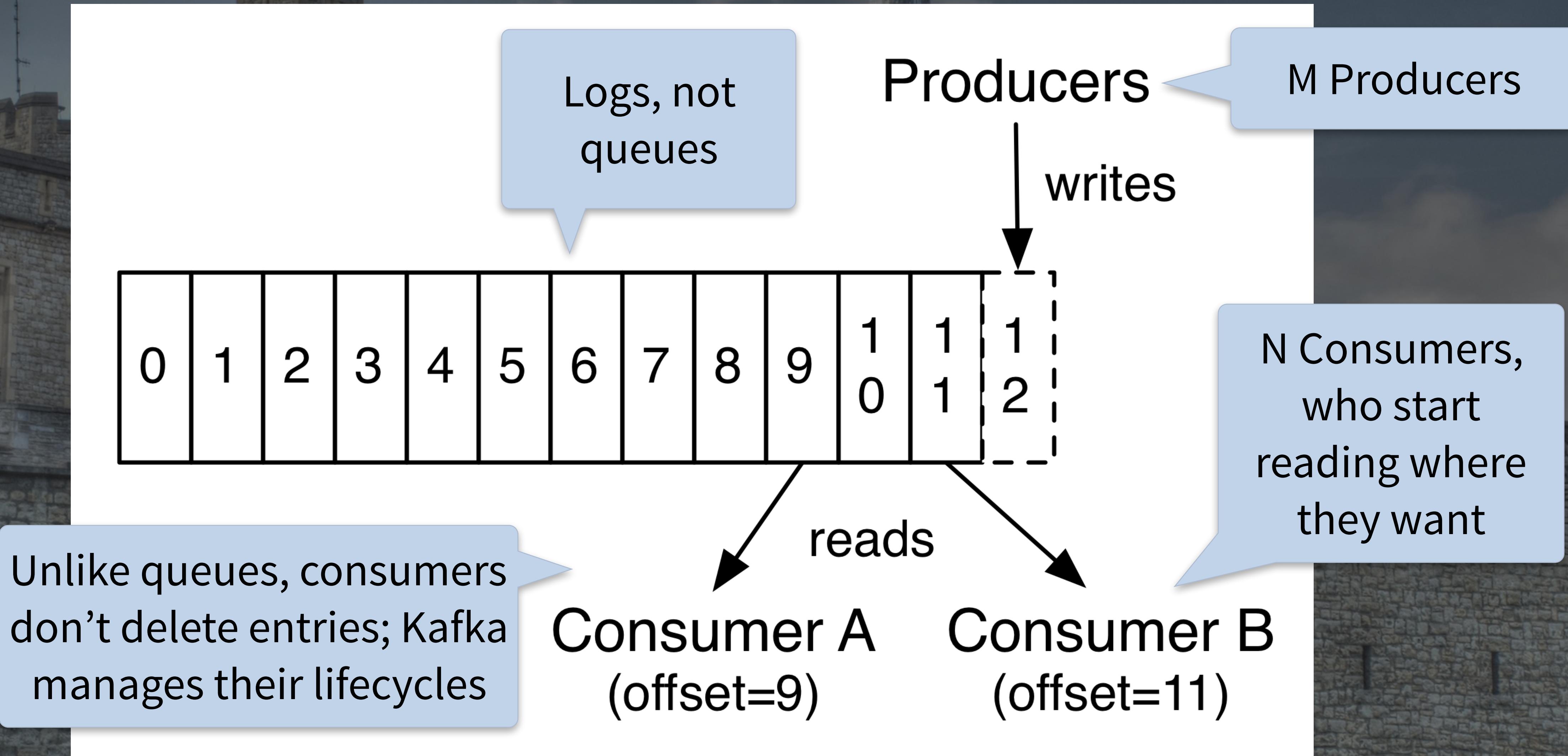


Topics are partitioned,
replicated, and
distributed

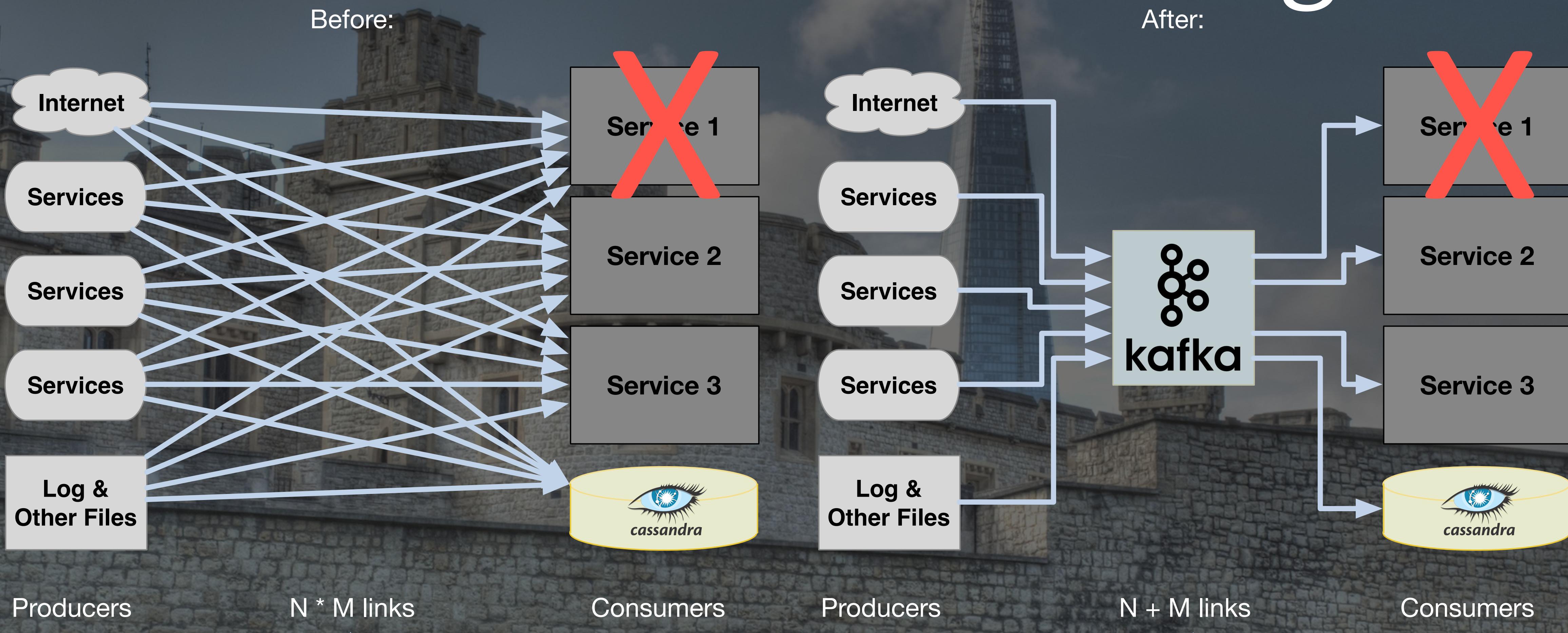
Consumer A
(offset=9)

Consumer B
(offset=11)

Why Kafka?

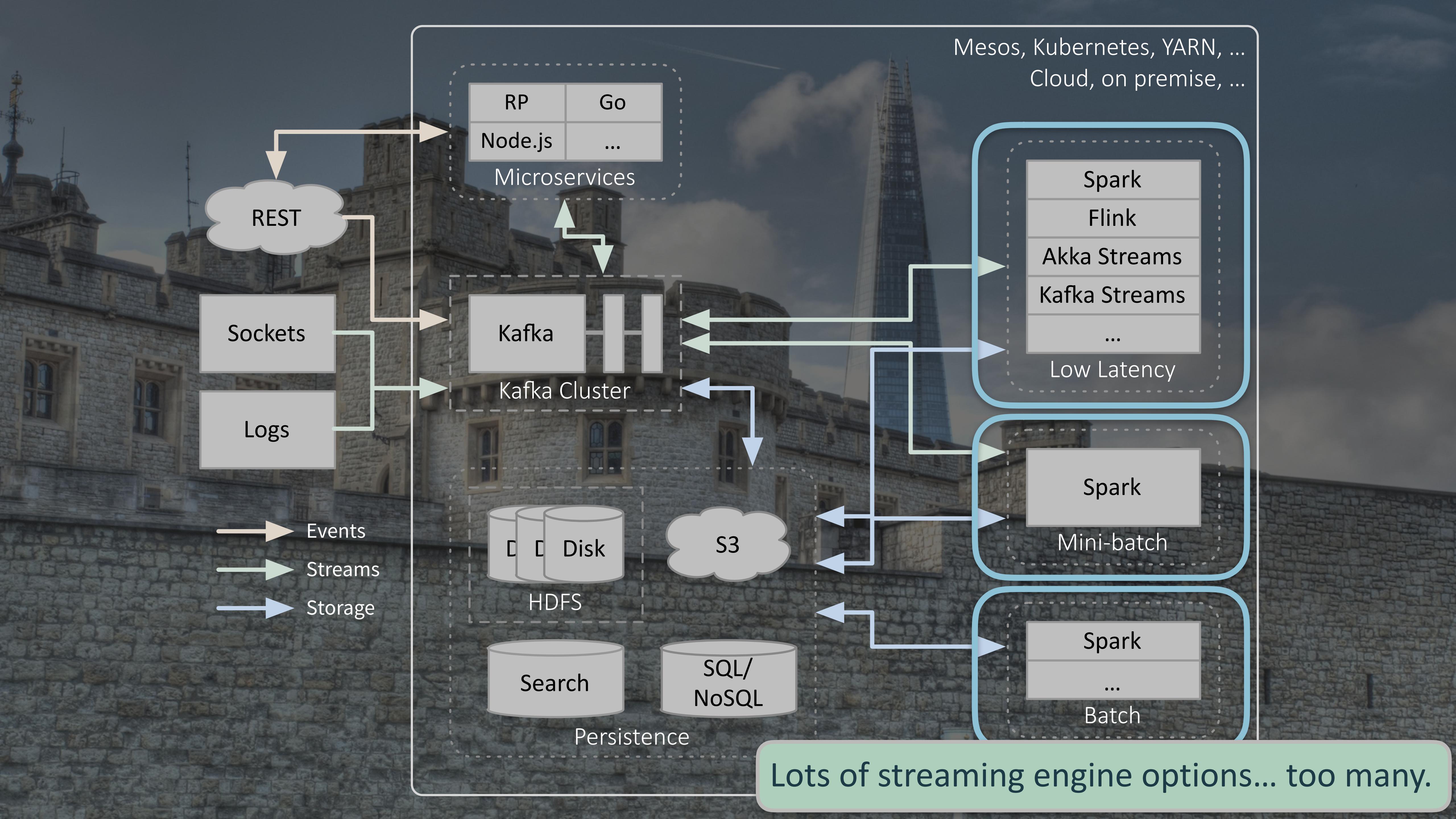


Using Kafka



Messy and fragile;
what if “Service 1”
goes down?

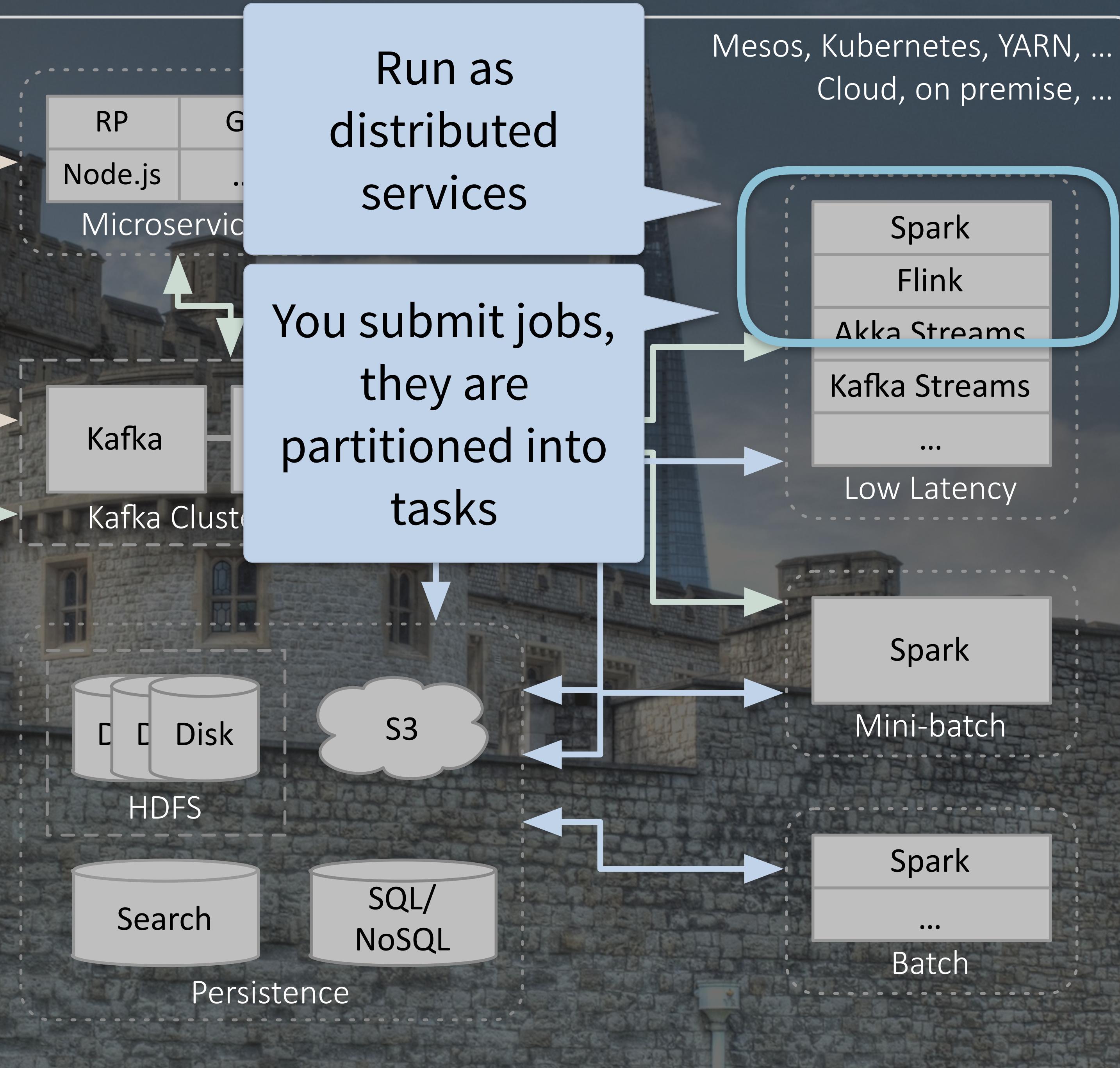
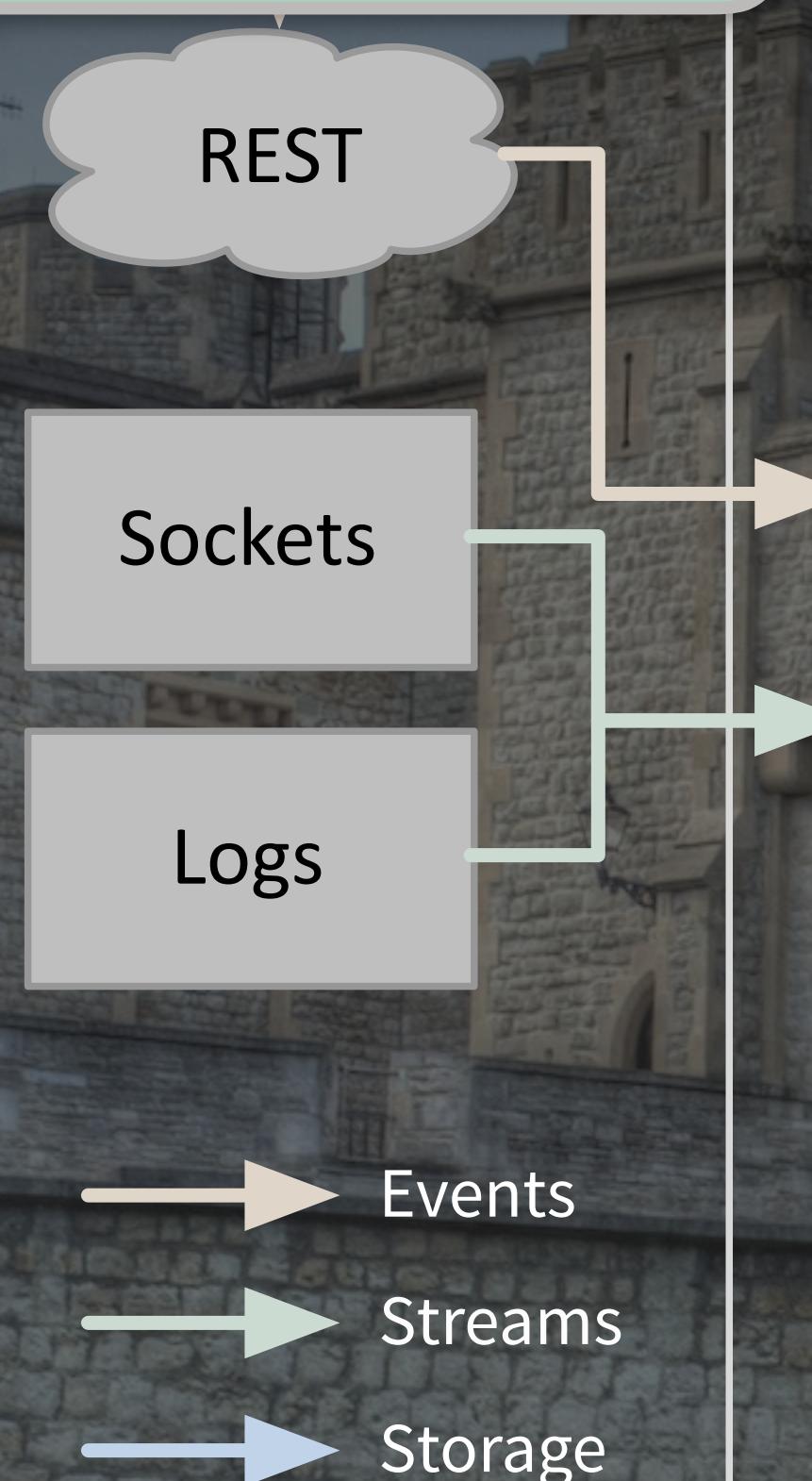
Simpler and more
robust! Loss of Service
1 means no data loss.



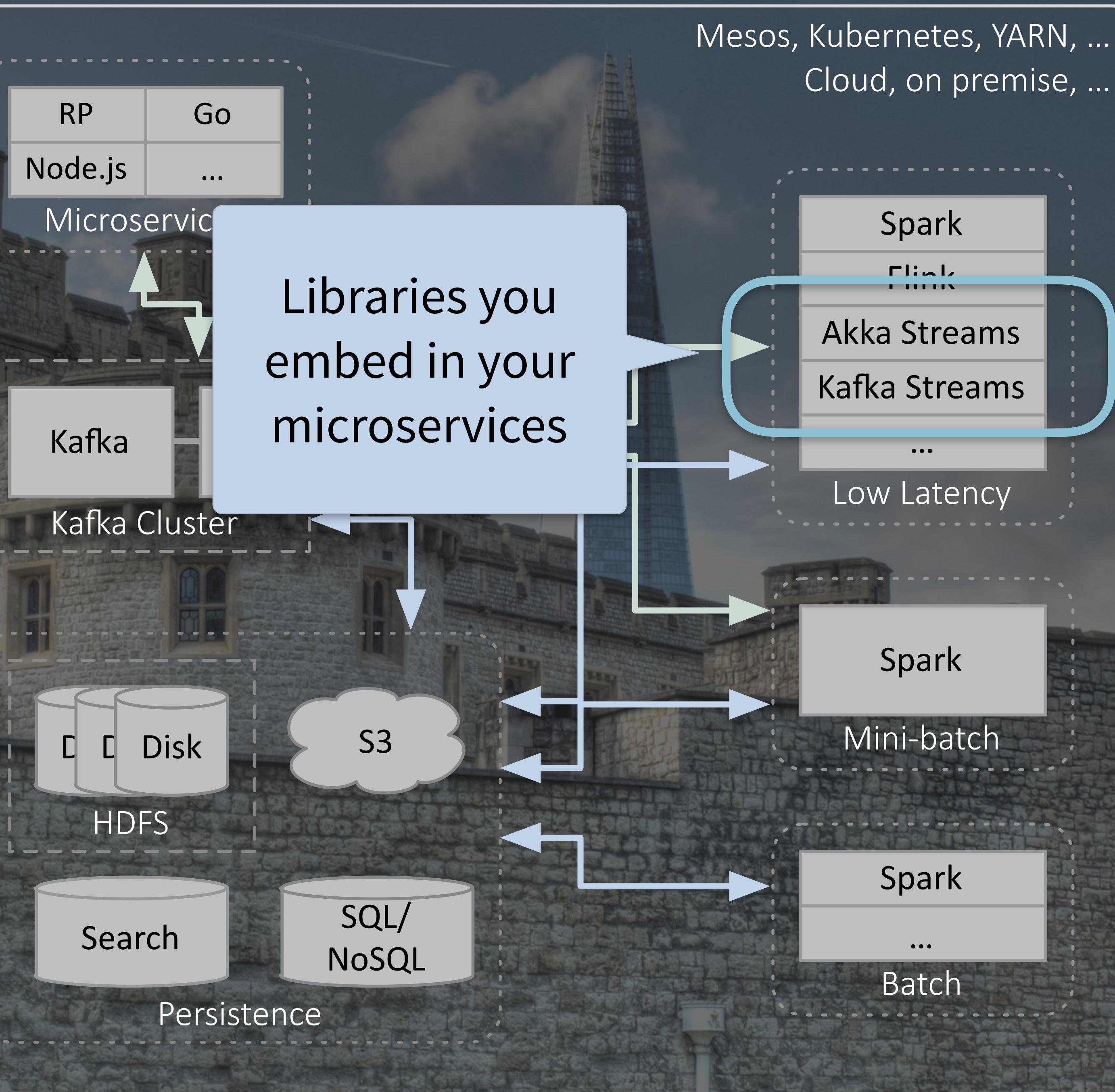
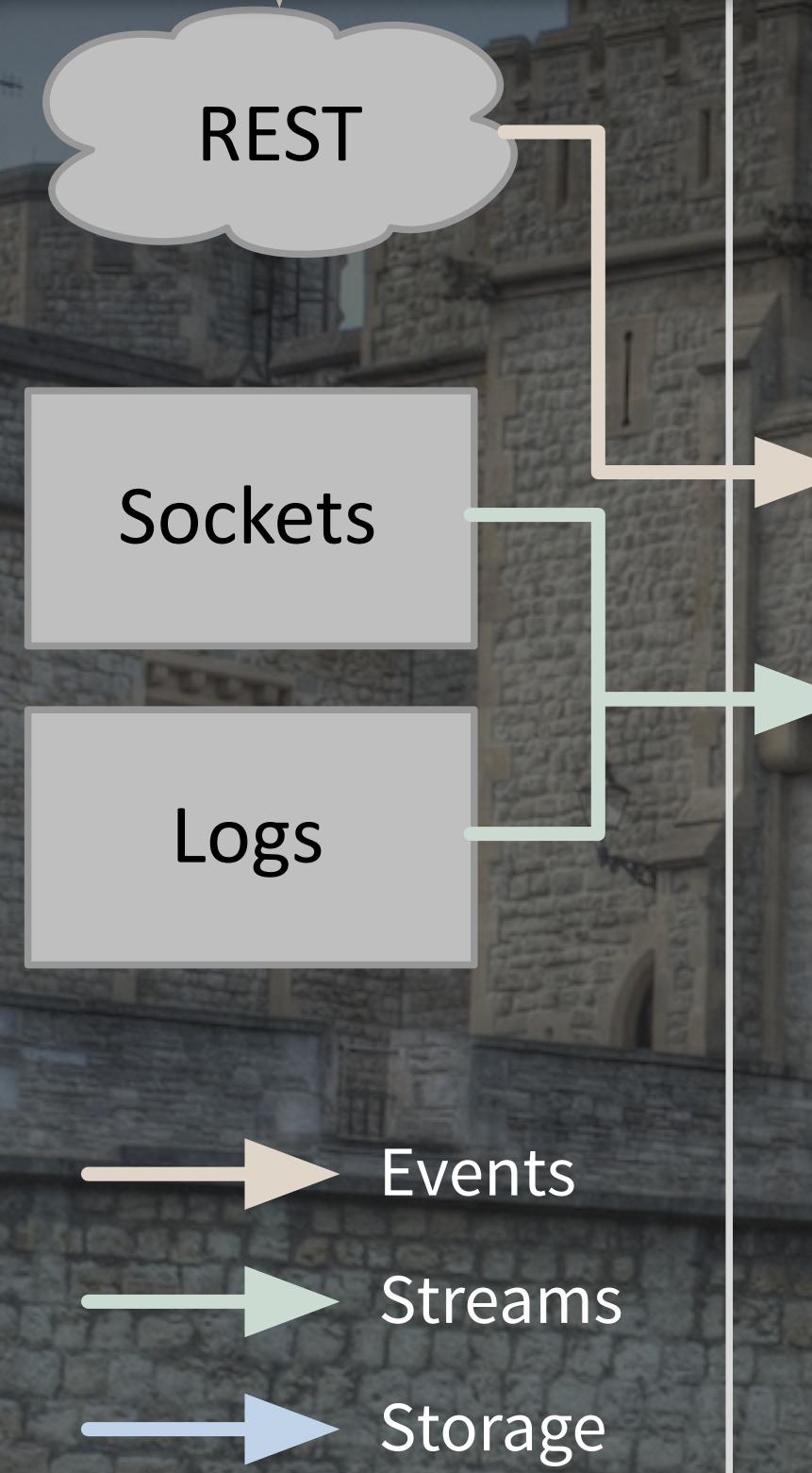
How do you choose?

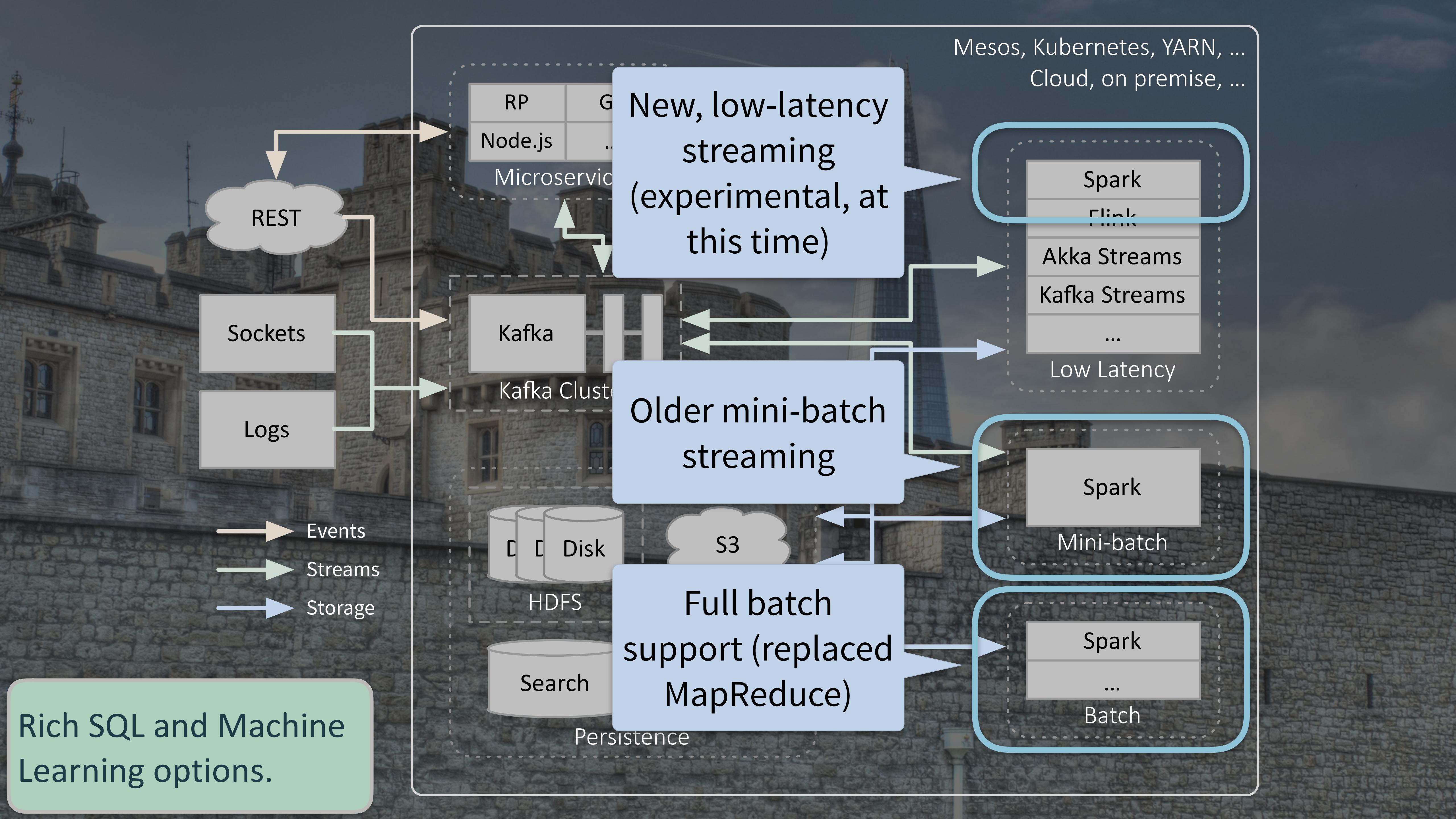
- Latency? How low?
- Volume: How high?
- Which kinds of data processing?
- How do you want to build, deploy, and manage these services?

The streaming engines form two groups:



The streaming engines form two groups:





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

New, low-latency
streaming
(experimental, at
this time)

Spark
Flink
Akka Streams
Kafka Streams
...

Low Latency

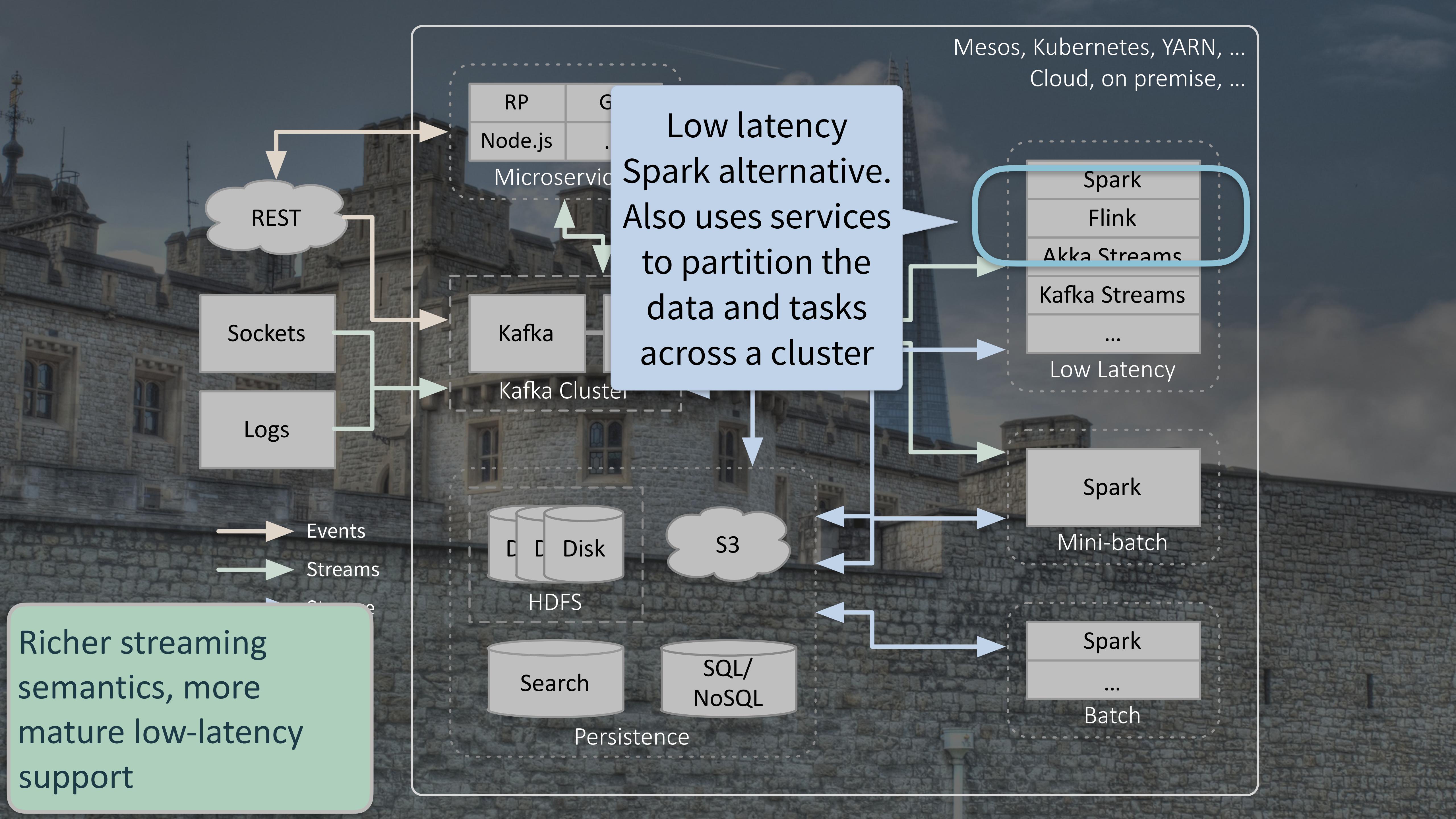
Older mini-batch
streaming

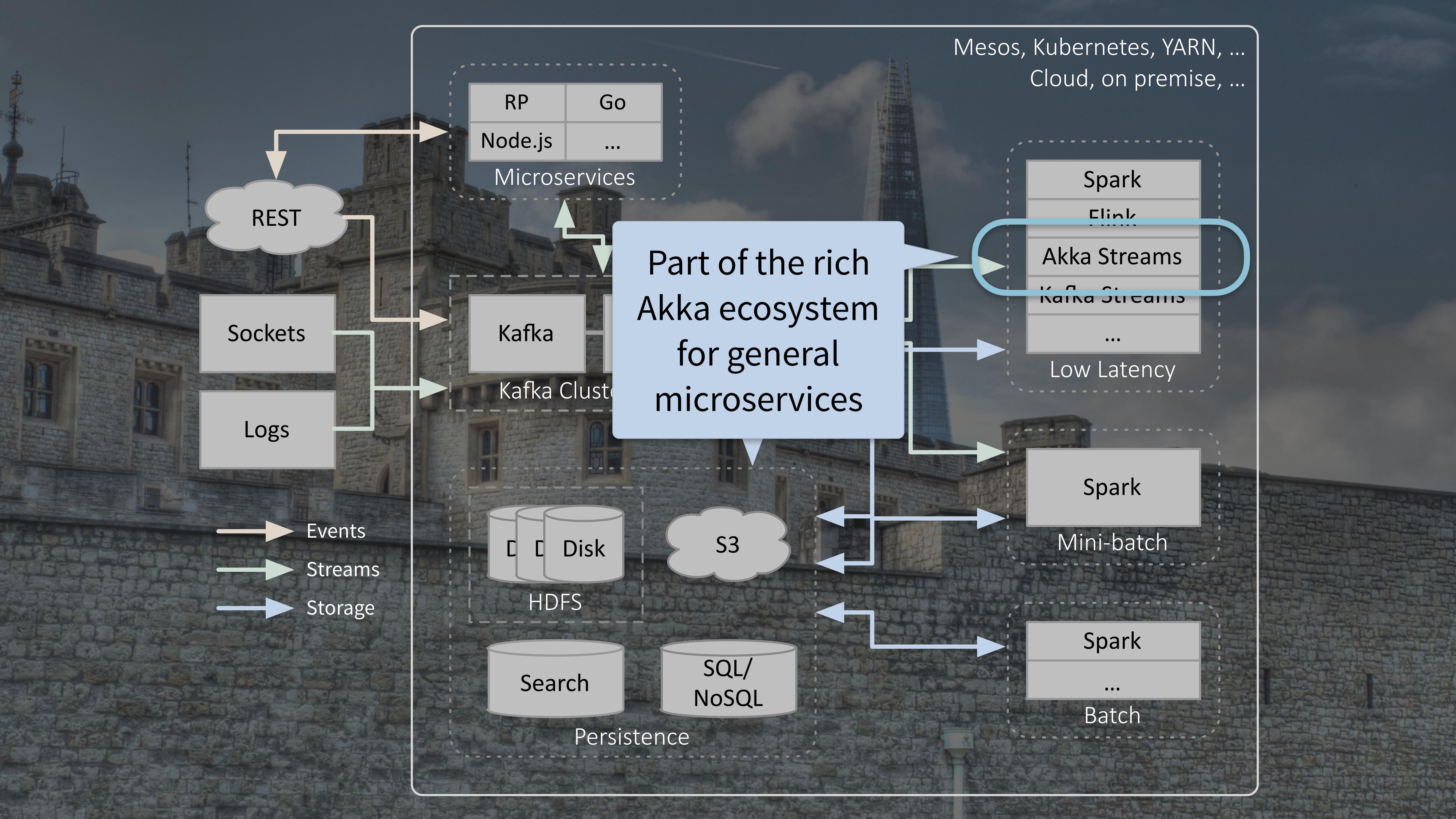
Spark
Mini-batch

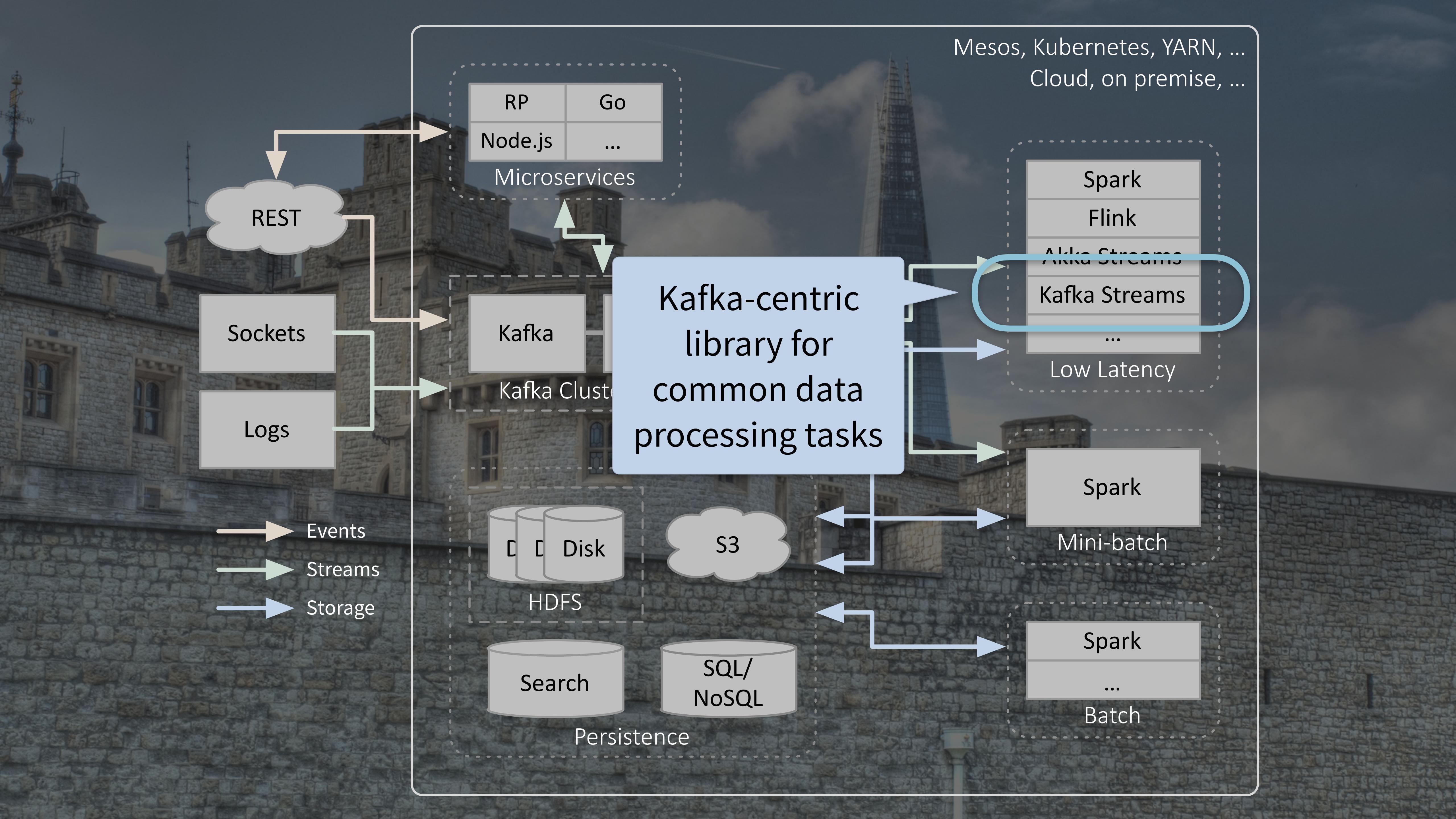
Full batch
support (replaced
MapReduce)

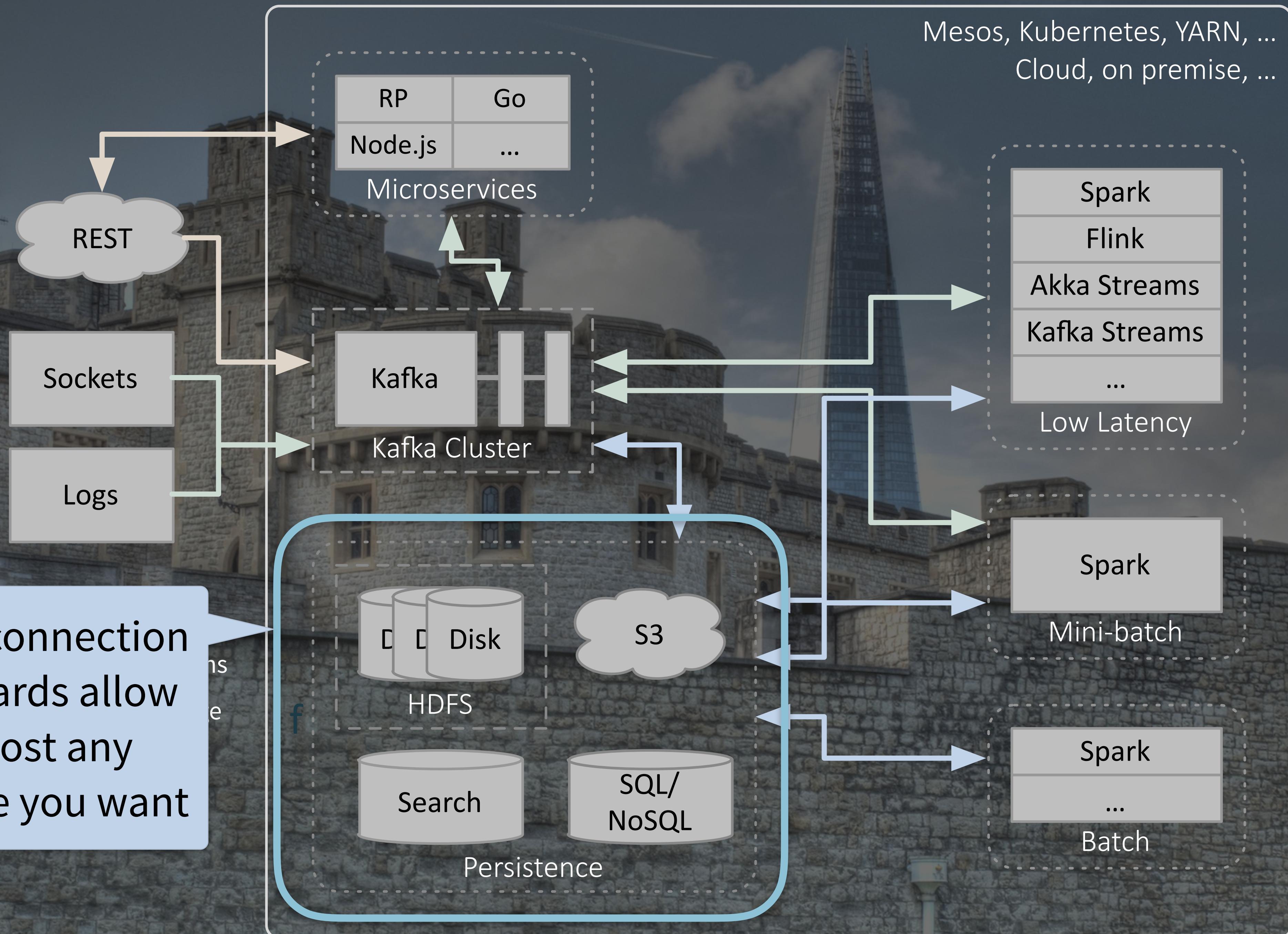
Spark
...
Batch

Rich SQL and Machine
Learning options.









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

Spark
Flink
Akka Streams
Kafka Streams
...

Low Latency

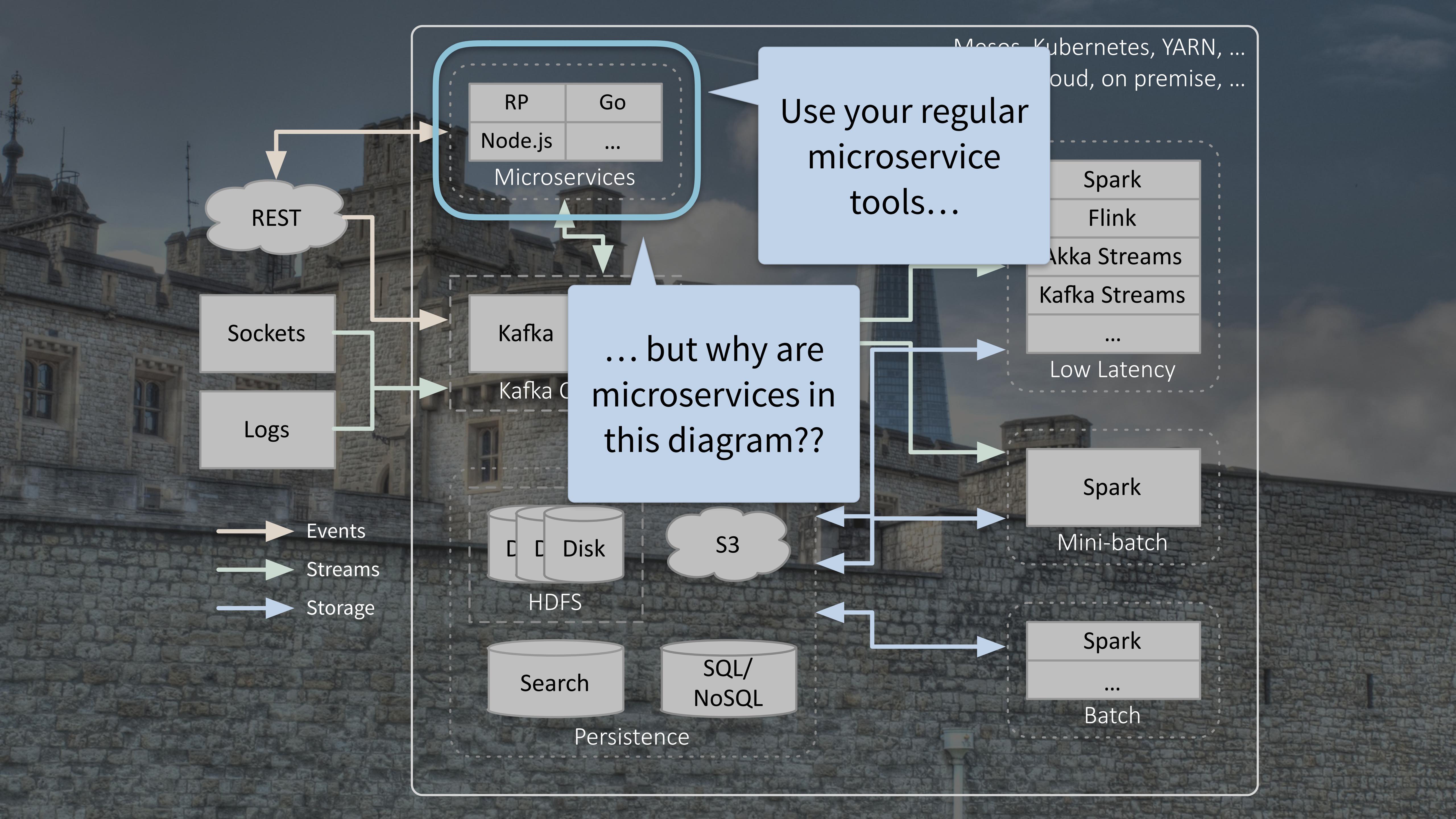
Spark

Mini-batch

Spark

Batch

Open connection
standards allow
almost any
storage you want



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?

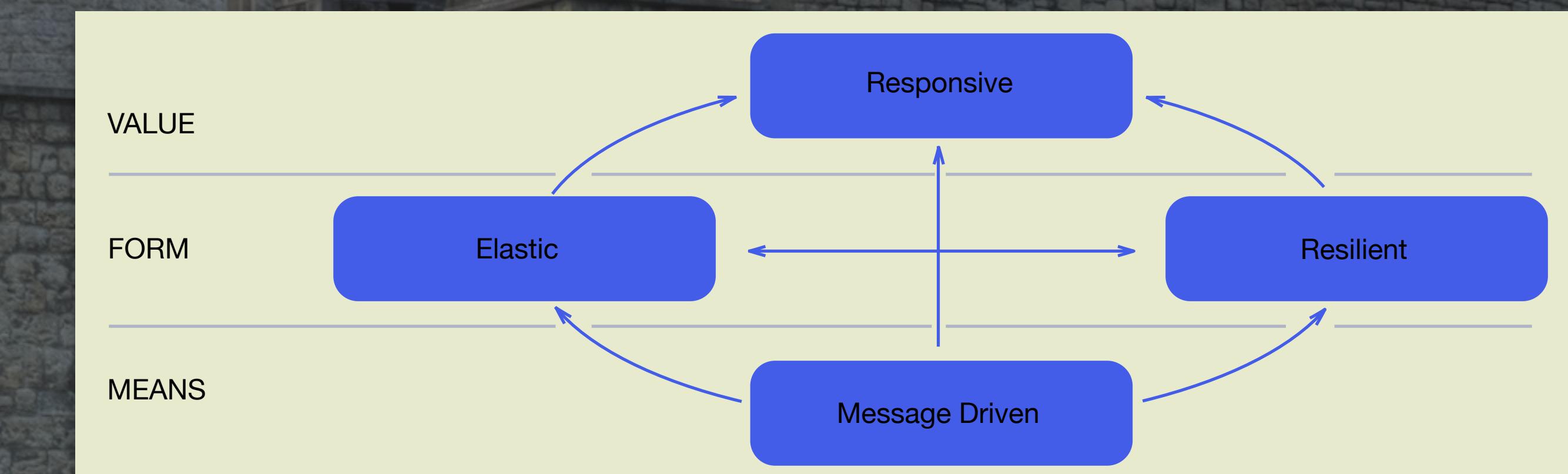
2. If streaming gives you information faster...

- ...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



Why Microservices in Fast Data?

4. This leads to our last major point...



Organizational Impact



Organizational Impact

- Data scientists have to understand deployment
- Data engineers have to become good at highly-available microservices
- Microservice engineers have to become good at data

The Past



Services

Big Data

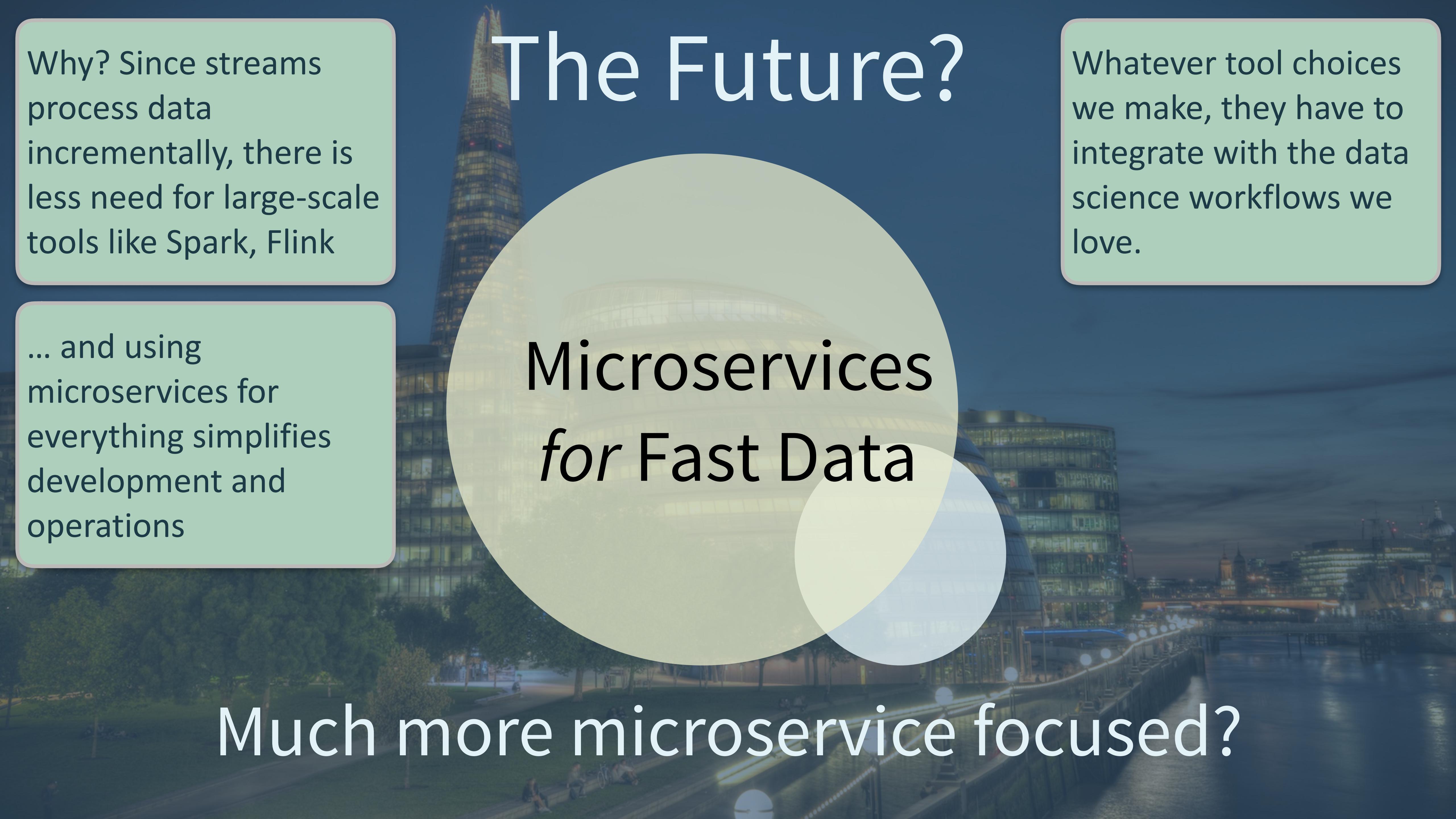
Some overlap in concerns, architecture

The Present



Microservices
& Fast Data

Much more overlap



Why? Since streams process data incrementally, there is less need for large-scale tools like Spark, Flink

... and using microservices for everything simplifies development and operations

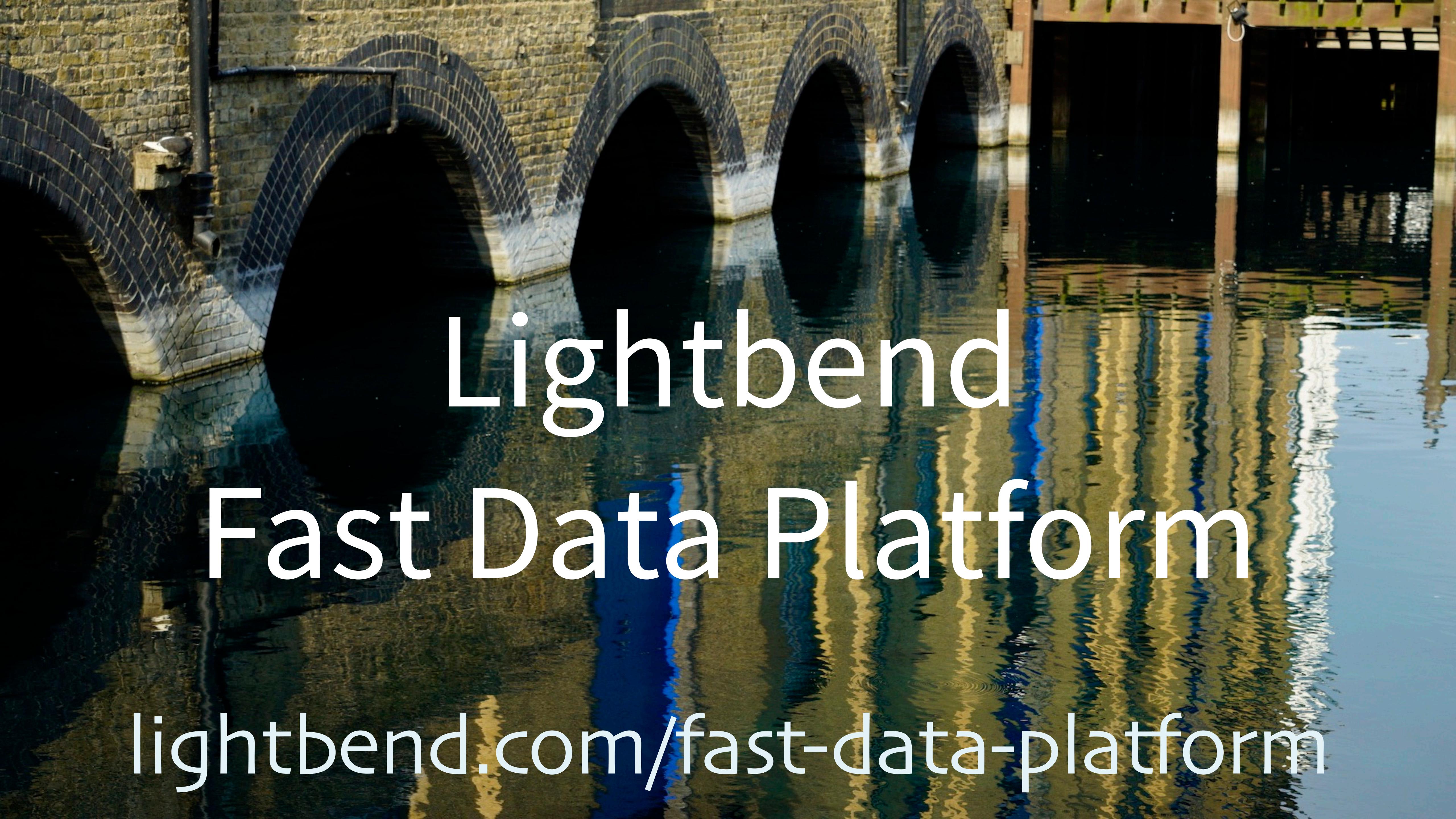
The Future?



Microservices
for Fast Data

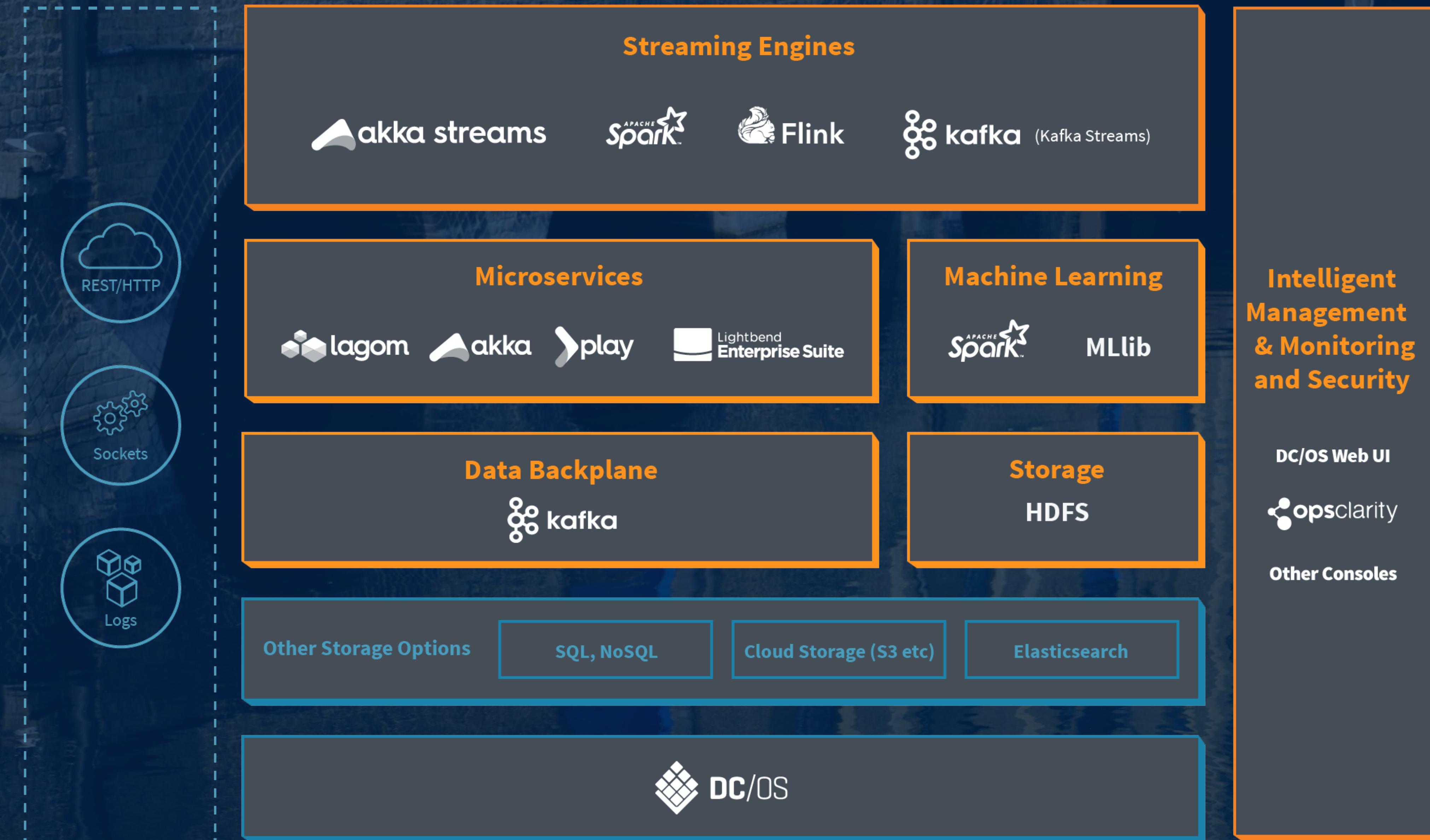
Much more microservice focused?

Whatever tool choices we make, they have to integrate with the data science workflows we love.

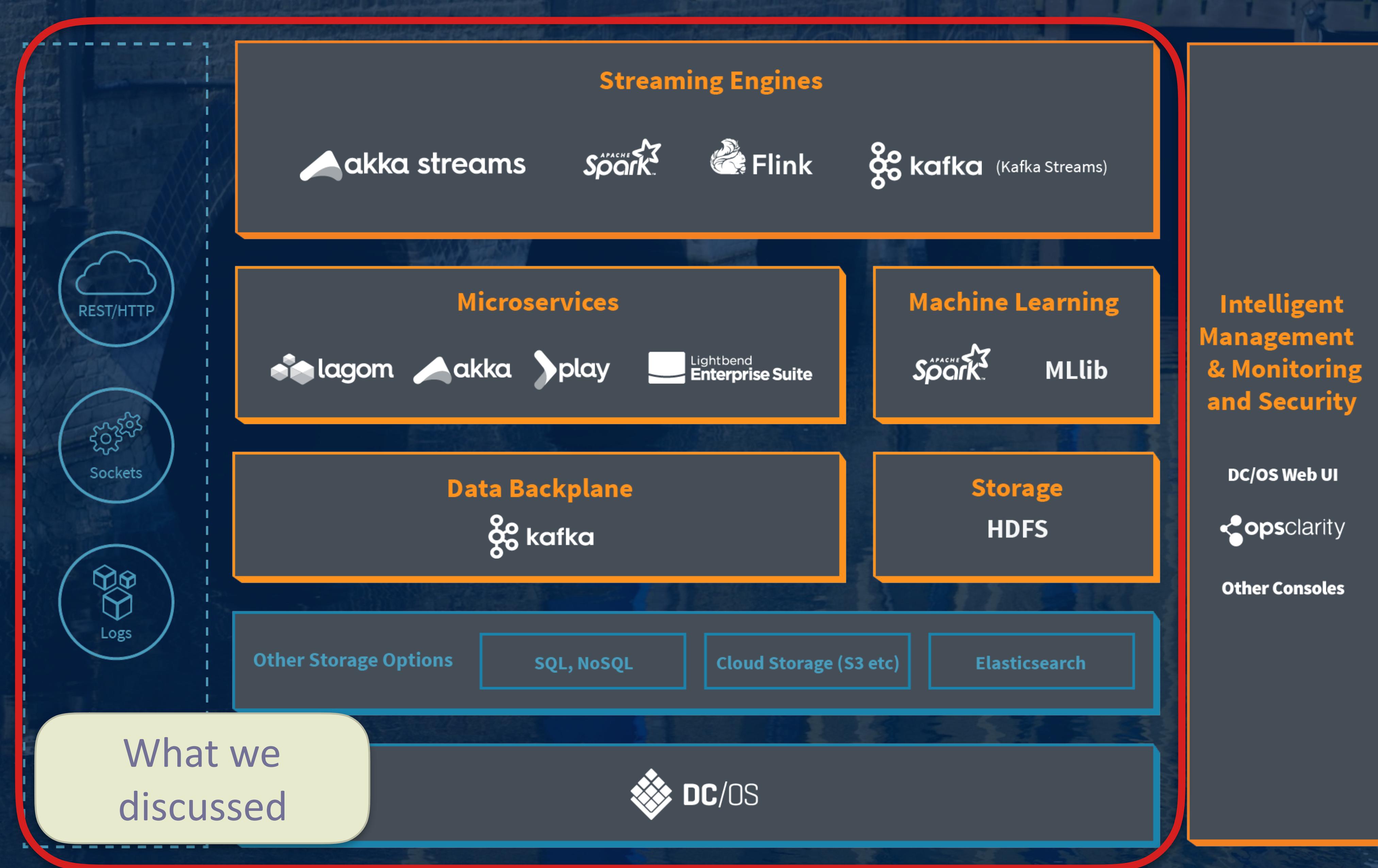
A photograph of a multi-arched stone bridge reflected perfectly in the still water below. The bridge's arches create a rhythmic pattern of light and shadow on the surface. A single bird is perched on a small ledge on the left side of the bridge. The overall scene is peaceful and symmetrical.

Lightbend Fast Data Platform

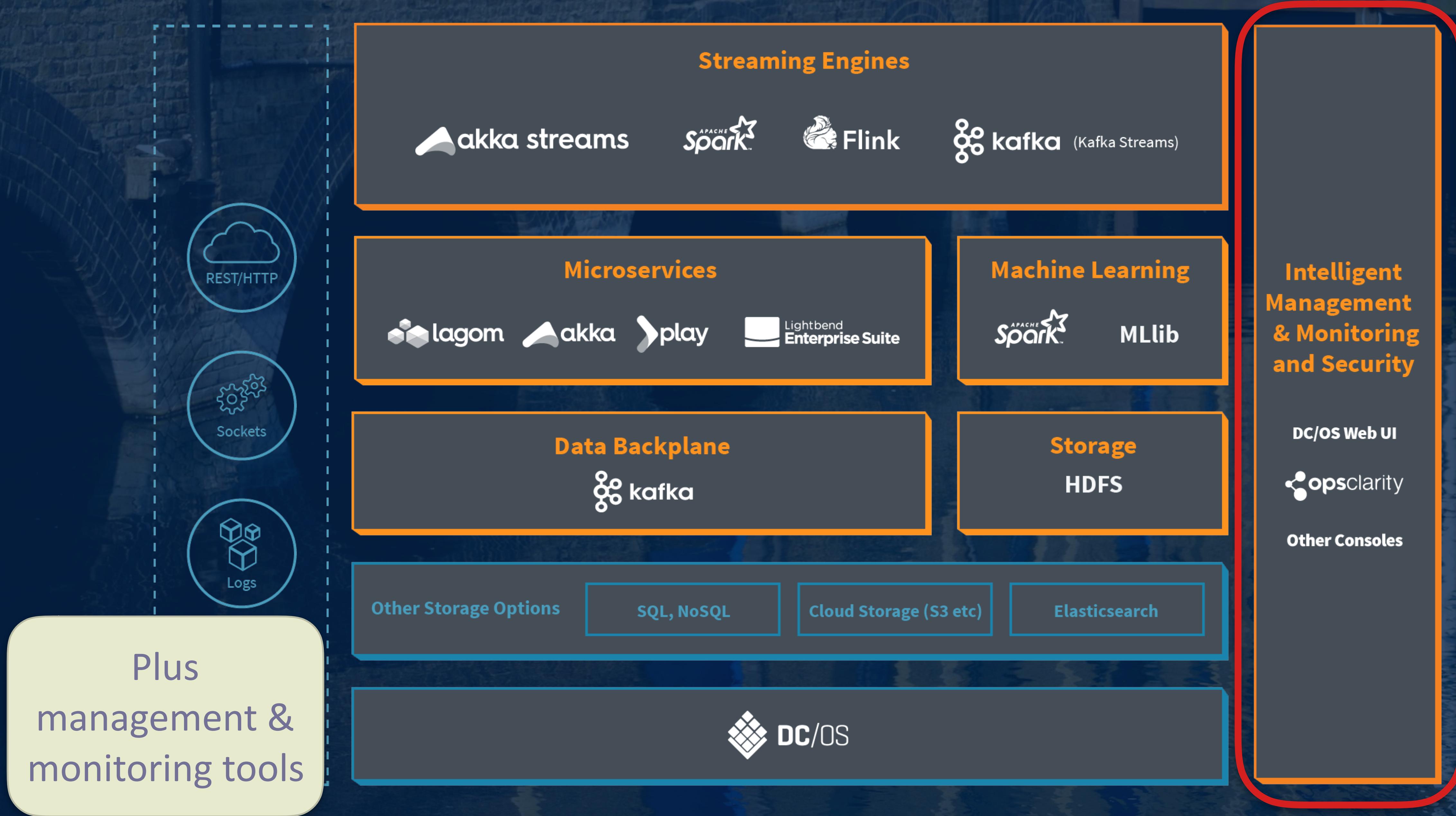
lightbend.com/fast-data-platform



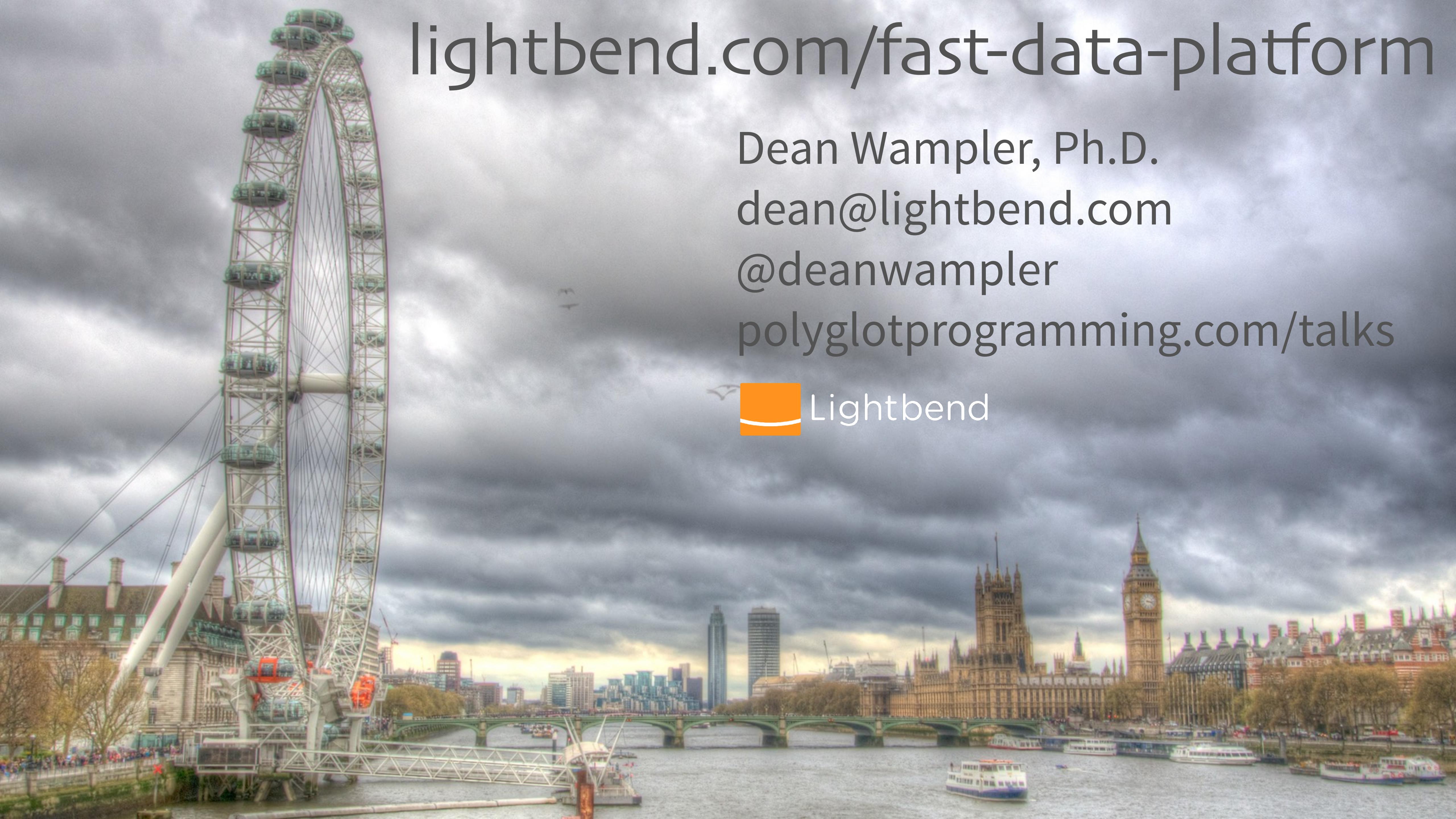
lightbend.com/fast-data-platform



lightbend.com/fast-data-platform



lightbend.com/fast-data-platform

A photograph of the London skyline featuring the London Eye (Millennium Wheel) on the left and the Palace of Westminster with Big Ben on the right. The River Thames is in the foreground with several boats. The sky is overcast with dramatic clouds.

lightbend.com/fast-data-platform

Dean Wampler, Ph.D.
dean@lightbend.com
[@deanwampler](https://twitter.com/deanwampler)
polyglotprogramming.com/talks

