

Executive Briefing: What it takes to use machine learning in fast data pipelines

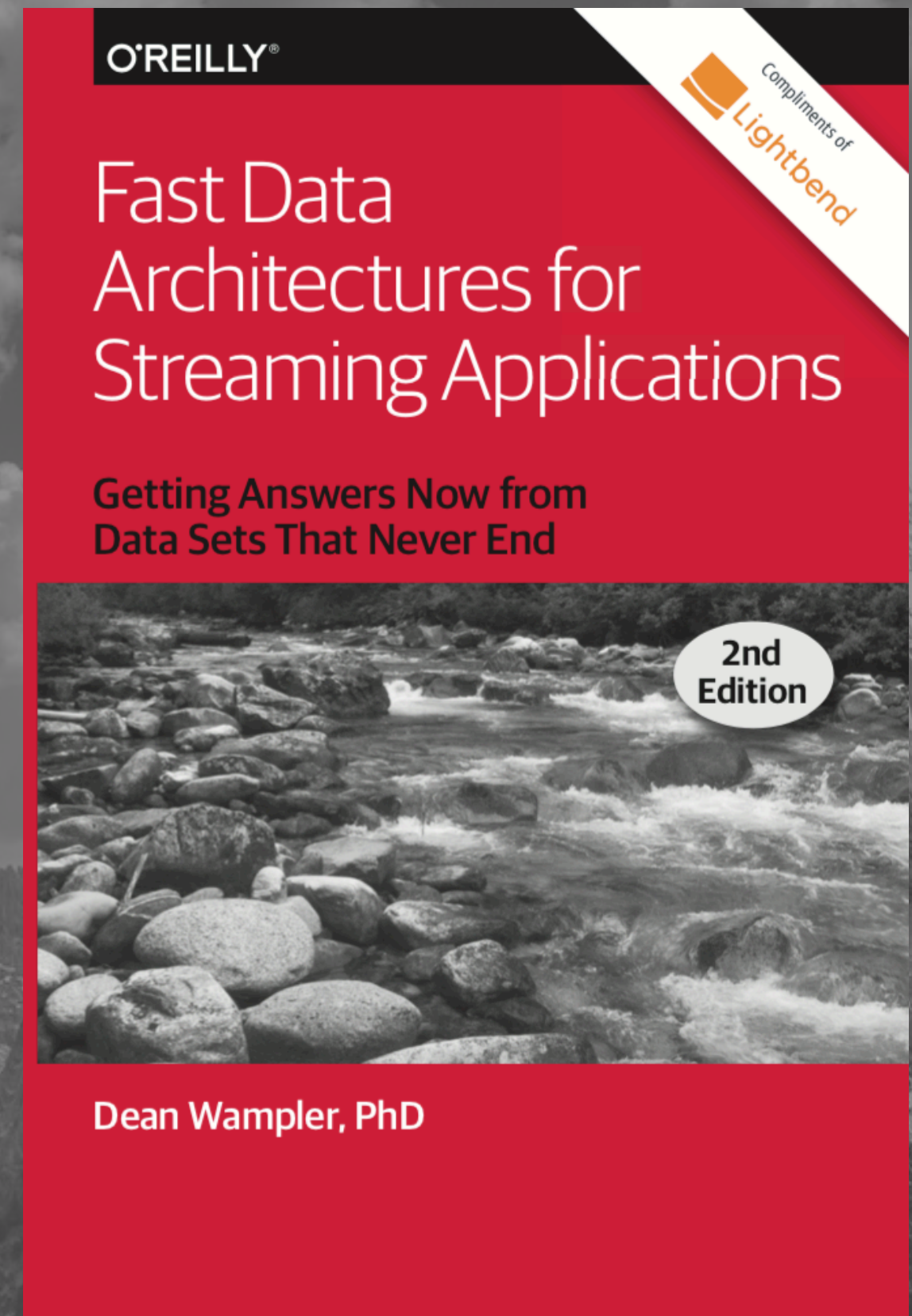
@deanwampler
dean@deanwampler.com
polyglotprogramming.com/talks

© Dean Wampler, 2007-2019

@deanwampler

Data Streaming, in General

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



@deanwampler

What We'll Discuss

- Batch vs. streaming... and why
- Data science vs. data engineering
- Serving models in production
 - CI/CD Systems for ML
 - Example architecture
 - Updating Models in Production

Batch vs. streaming... and why



Telecom



Finance



Medical

Energy

... and IoT

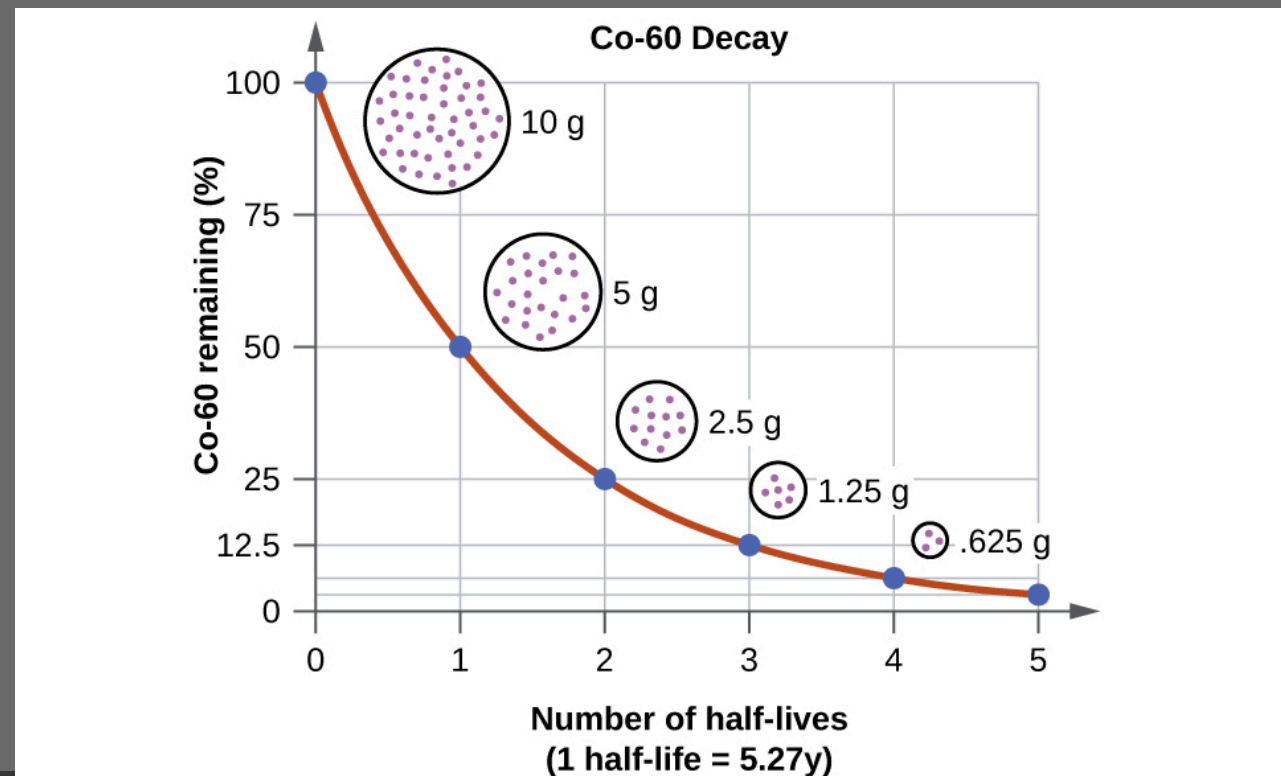
State of the art phone!




Mobile

@deanwampler

Information value
has a half life;
it decays with time

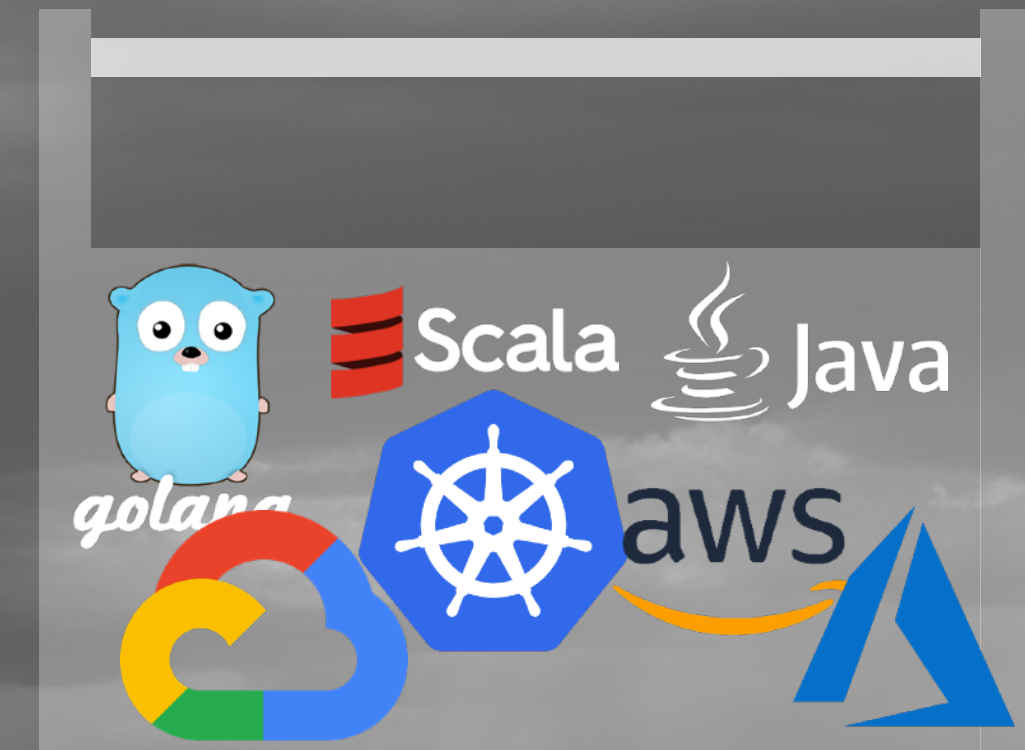




Data Science vs. Data Engineering



Data Science toolbox



Software Engineering toolbox

Data Scientists

- Comfortable with uncertainty
- Less process oriented
 - Iterative, experimental



Data Engineers

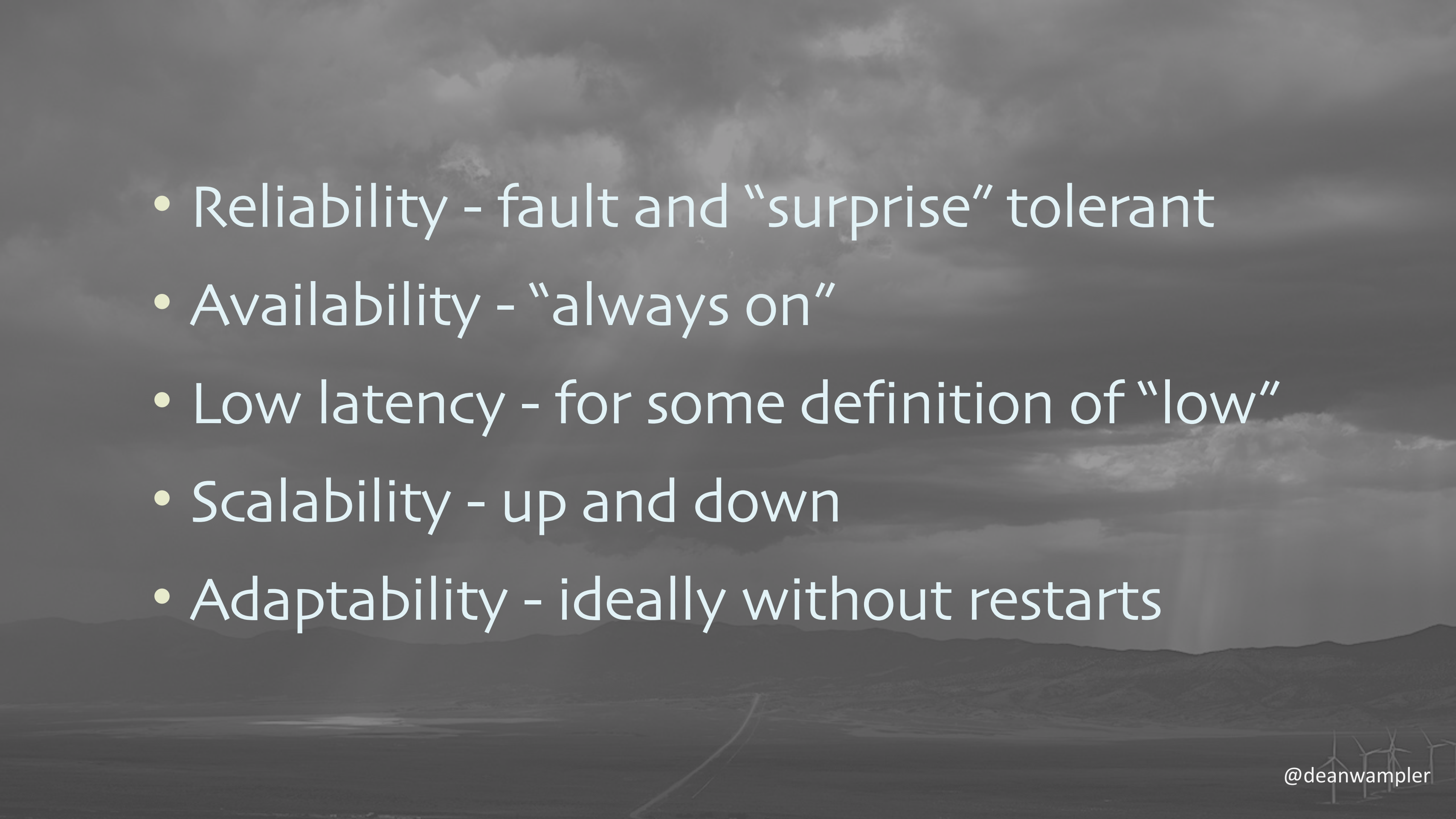
- Uncomfortable with uncertainty
- Process oriented
 - Agile Manifesto
 - ... which does not mention data!



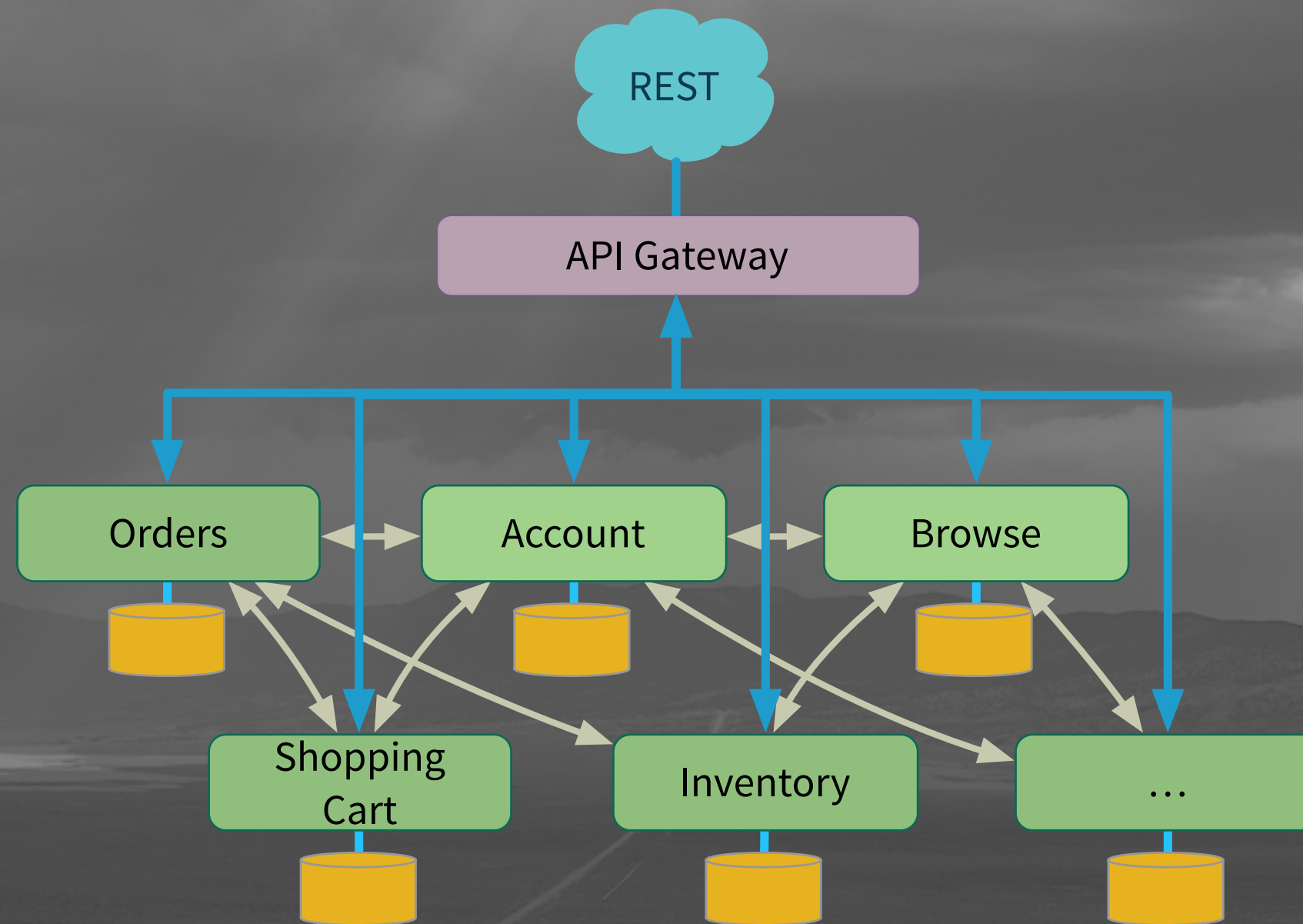
<https://derwen.ai/s/6fqt>

Streaming Imposes New Requirements

If you run something long enough, all rare problems eventually happen!

- 
- Reliability - fault and “surprise” tolerant
 - Availability - “always on”
 - Low latency - for some definition of “low”
 - Scalability - up and down
 - Adaptability - ideally without restarts

In other words: Microservices



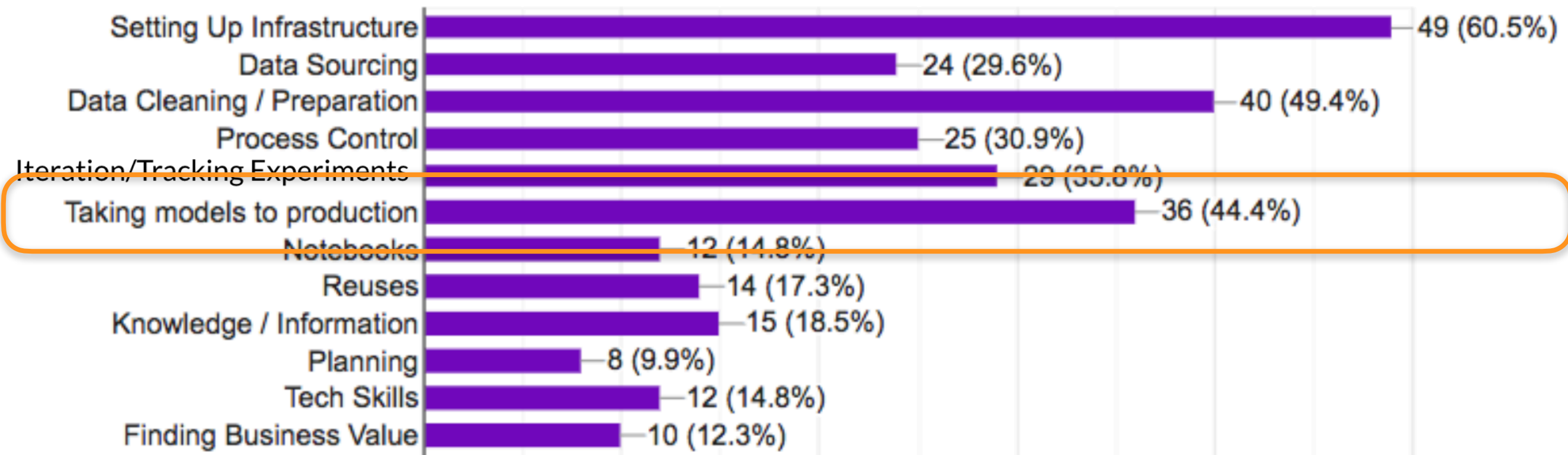
Serving Models in Production

A Recent Kubeflow User Survey

What are the major pain points in your ML workflows today? Check all that apply.



81 responses



Kubeflow User Survey - 1Q2019

From this [Kubeflow Overview](#)

@deanwampler

Lack of Tool/Process Integration

- ~60% worry about missed opportunities
- ~50% worry about loss of data team productivity
- ~45% worry about slow time-to-market
- ~40% worry about customer dissatisfaction

Can You Answer this Question?

- Why did the model reject that loan application?
(After you've been sued for discrimination...)

Which model was it?

- Which **version** of the model was used?
- How was it **trained**?
- When was this model **deployed**?
- ... and other questions you'll need to answer to understand what happened...

A close-up photograph of a broken feather lying on a dark, pebbly ground. The feather is split into two main parts: a larger, brownish, frayed section on the left and a smaller, white, clean section on the right. The ground is composed of small, dark, rounded pebbles and some dry organic matter. The lighting is bright, casting shadows and highlighting the textures of the feather and the ground.

CI/CD for ML?

a.k.a “MLops”

@deanwampler

CI/CD Process Required (1/4)

- Version control - for models and code
- Automation - builds, tests, quality checks, artifact management & delivery
 - Necessary for reproducibility

CI/CD Process Required (2/4)

- Supports different launch configurations:
 - “dark” launches
 - A/B, Canary, and other testing scenarios

CI/CD Processes Required (3/4)

- Auditing
 - Which model used to score this record?
 - Which records used to train this model?
 - Who accessed this model and when?

CI/CD Processes Required (3/4)

- Auditing
 - Which model used to score this record?
 - Which records used to train this model?
 - Who accessed this model and when?

Models Are
Data

CI/CD Processes Required (3/4)

- Auditing
 - Which model used to score this record?
 - Which records used to train this model?
 - Who accessed this model and when?

GDPR - What if a customer asks you to delete their data? Do you also delete the models trained with that data?

CI/CD Processes Required (4/4)

- Monitoring
 - Resource utilization changes?
 - Quality metrics:
 - Match performance during training?
 - Concept drift?

What's Different from Microservice CI/CD?

- AutoML
- Data safety and lineage
- Model fairness and reproducibility
- Model and feature artifact management

<https://www.oreilly.com/ideas/9-ai-trends-on-our-radar>

@deanwampler

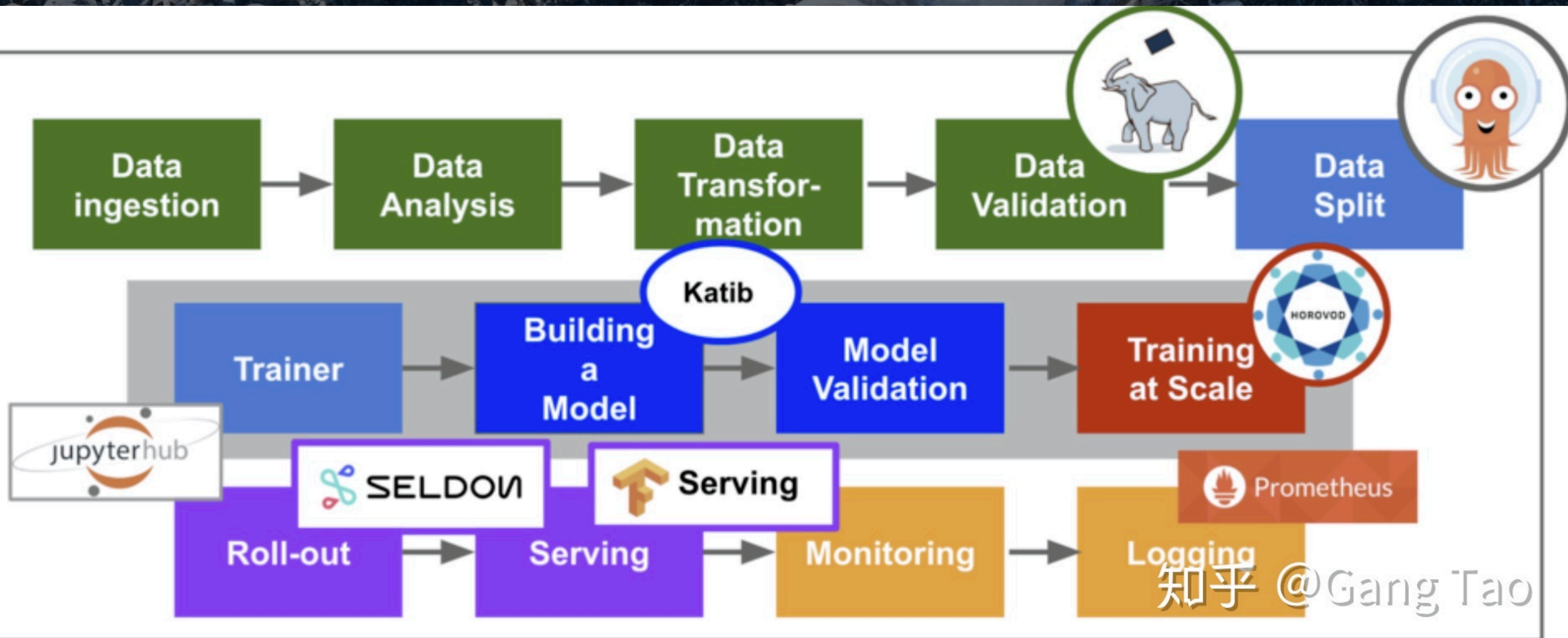
What's Different from Microservice CI/CD?

- CI/CD prefers deterministic measures of quality. How should you support the extra statistical indeterminacy data science introduces?

CI/CD Suites for ML

- Kubeflow - for Kubernetes
- SageMaker - for AWS users
- MLFlow - from the Spark community
- ... plus emerging vendors

Systems - Kubeflow



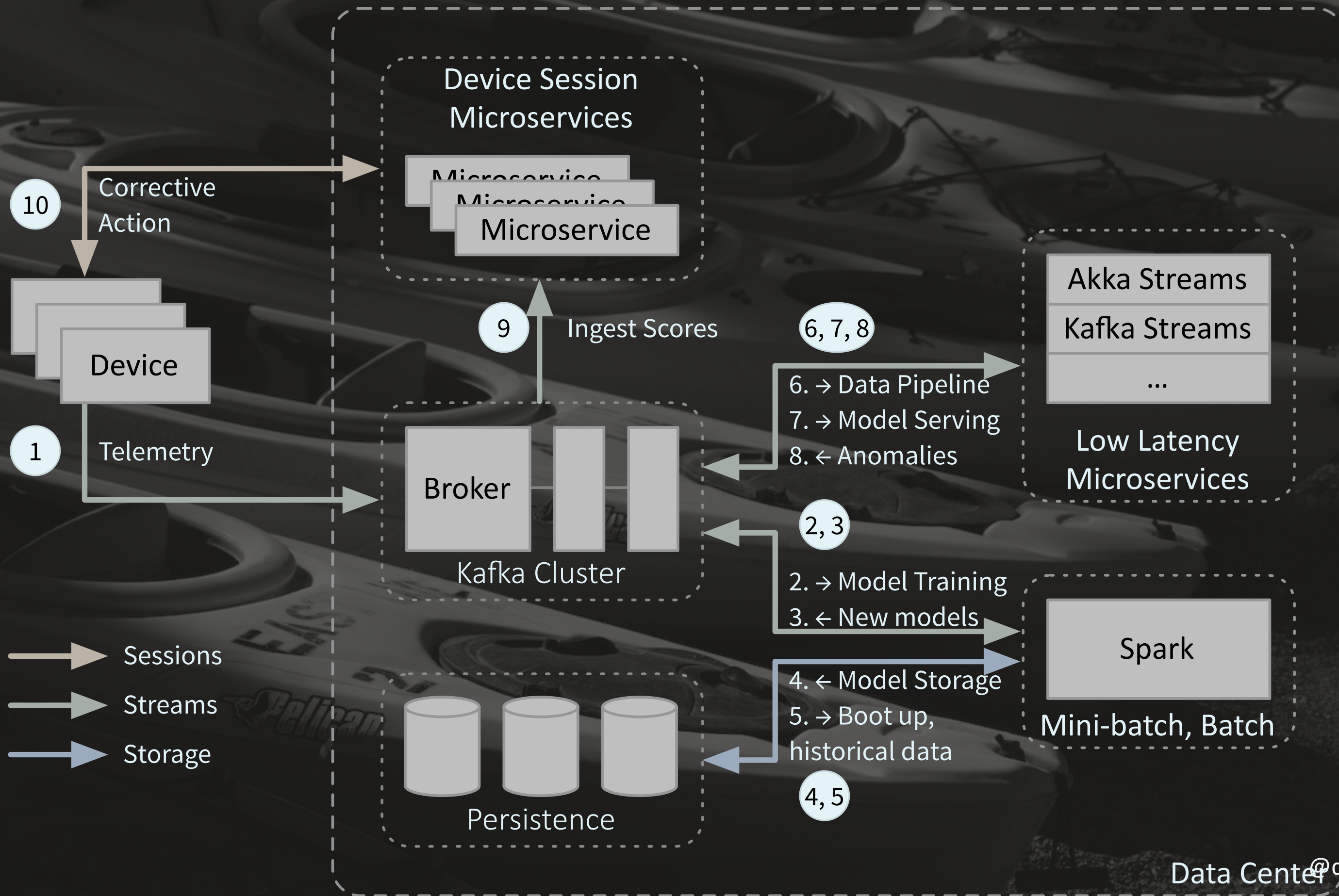
Example Architectures

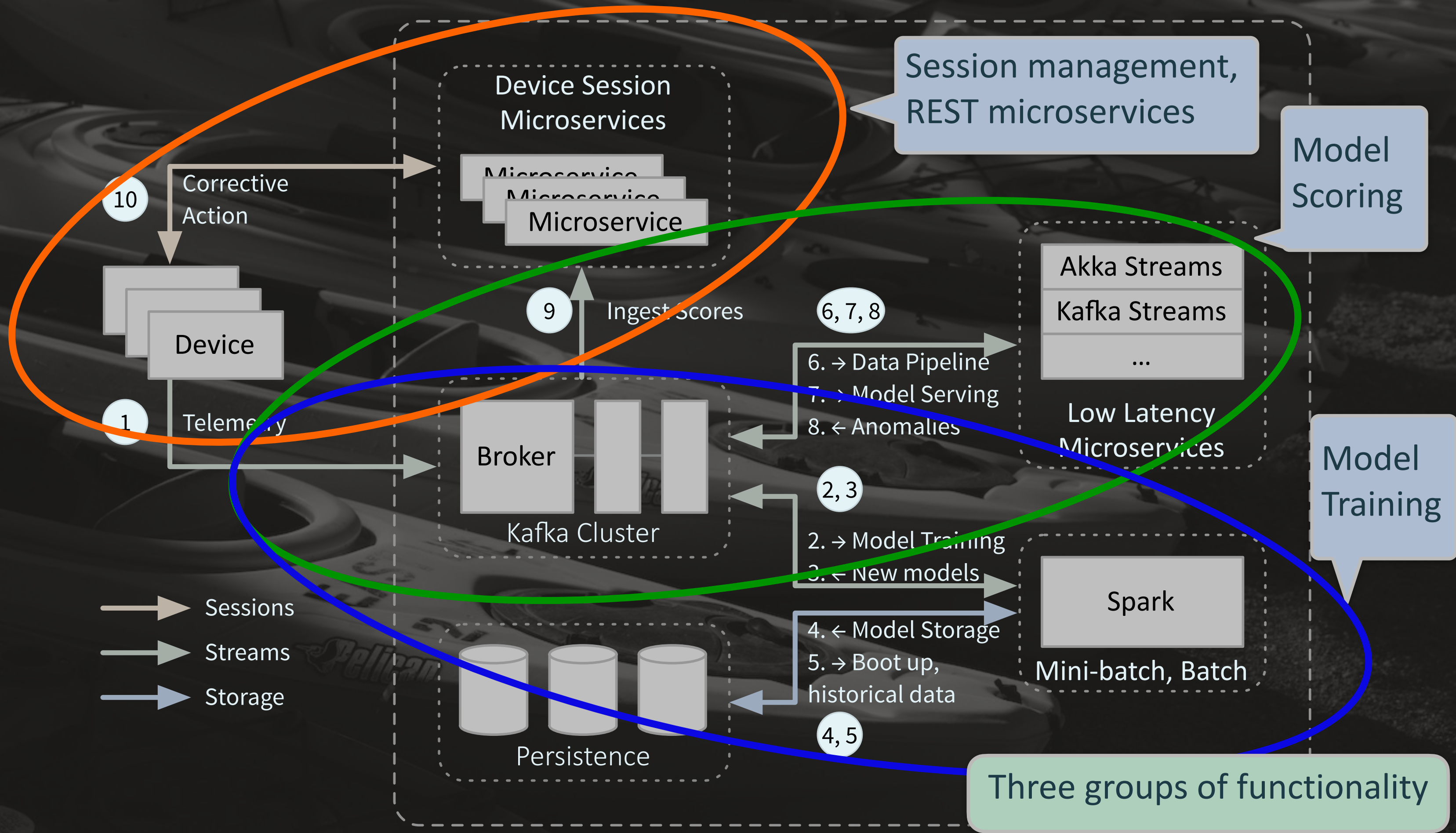


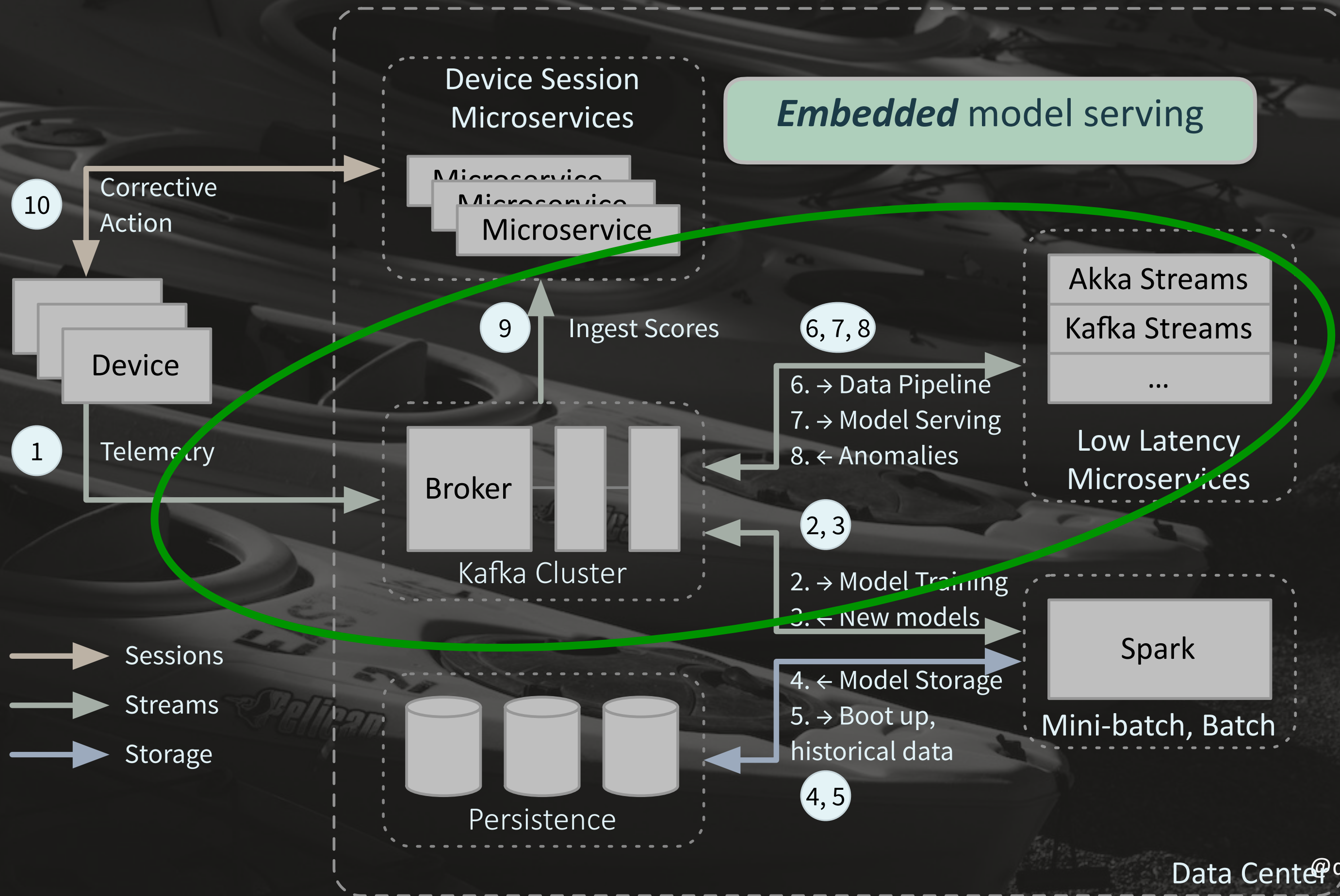
A black and white photograph of several white Pelican kayaks lined up on a sandy beach. The kayaks are viewed from a high angle, showing their hulls and cockpit openings. The brand name 'Pelican' is visible on the side of the kayaks in the foreground. The background shows more kayaks receding into the distance.

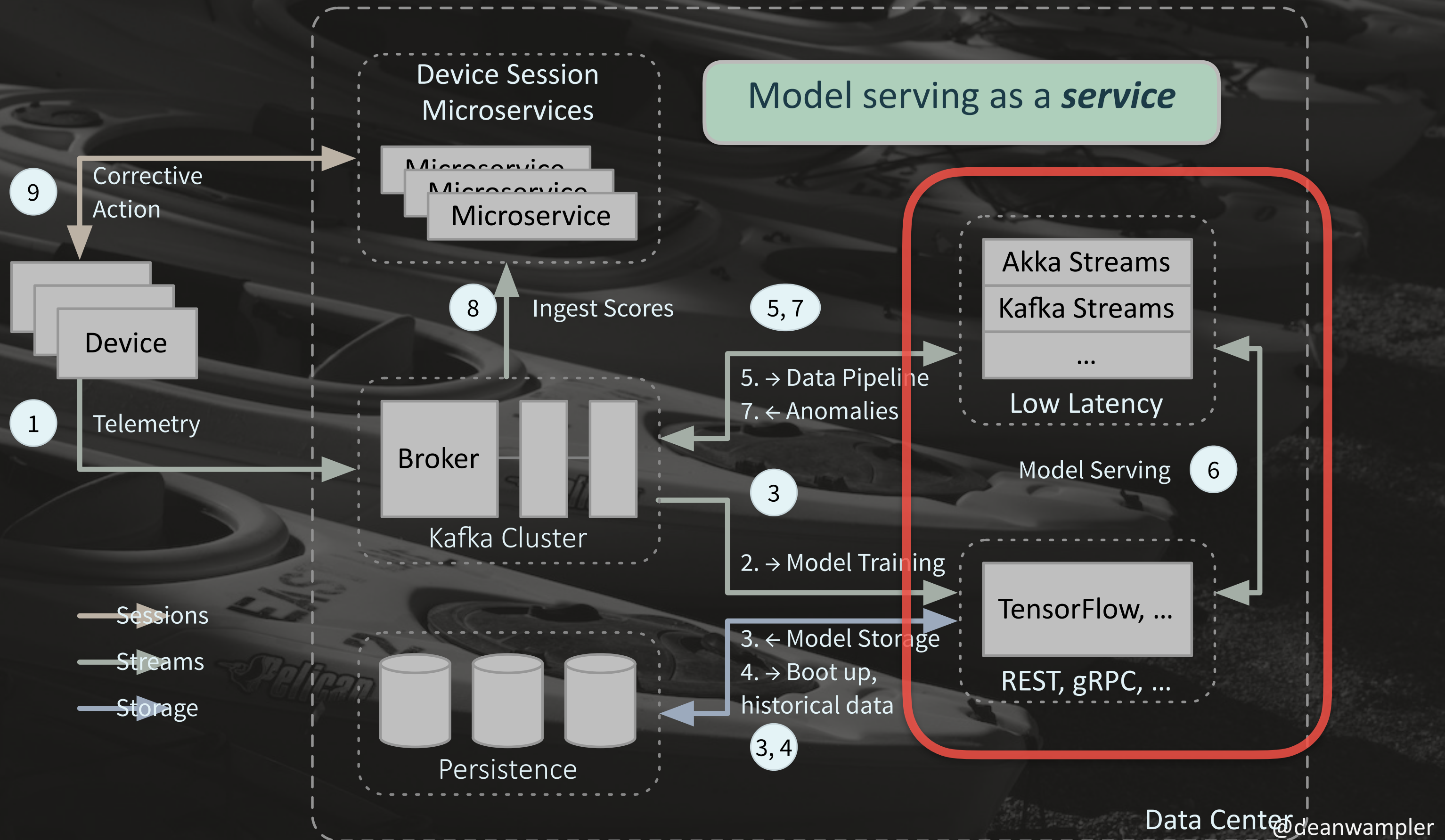
Example Architectures

Timely Information
Integrated with
Your Apps

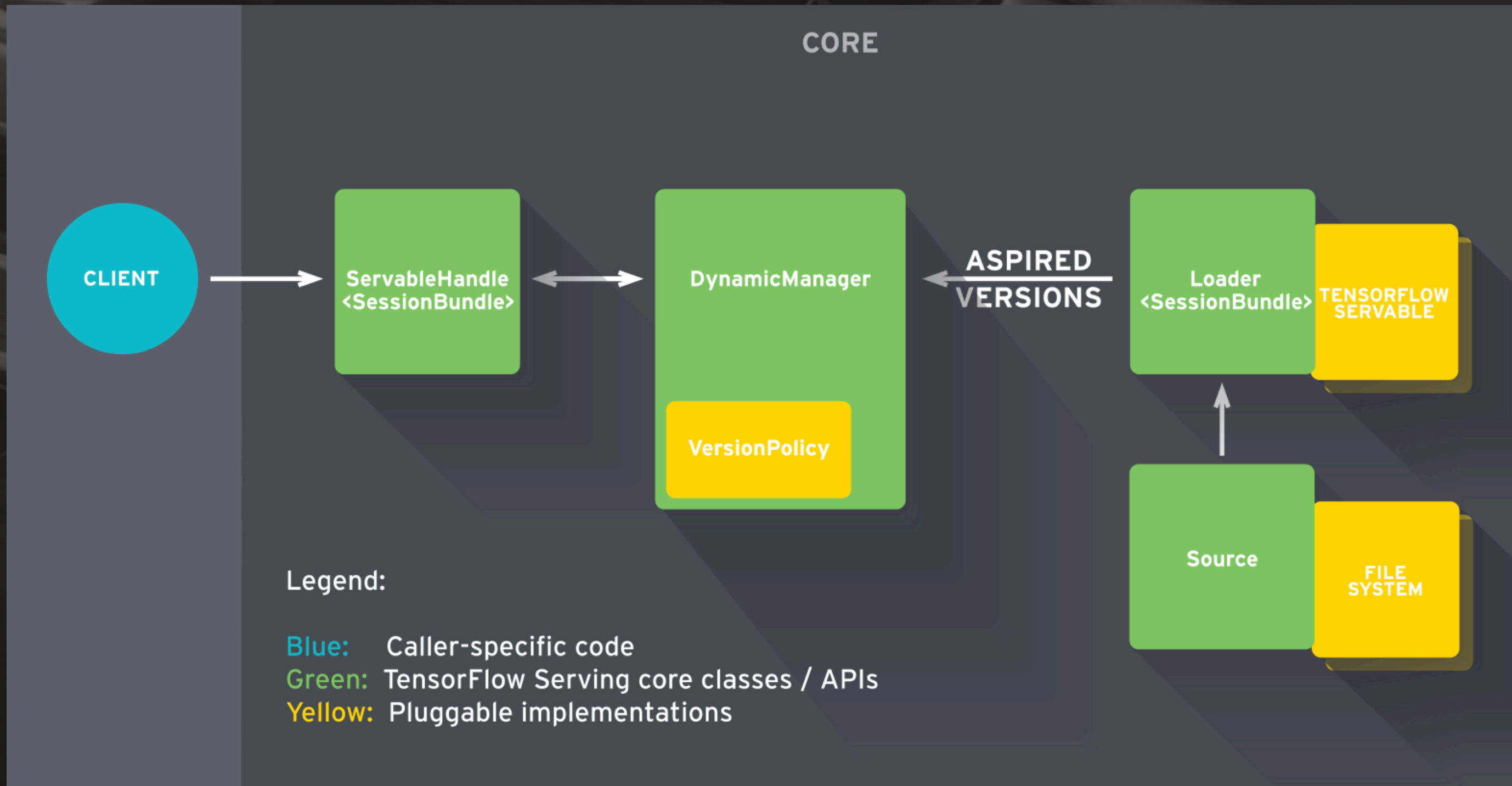






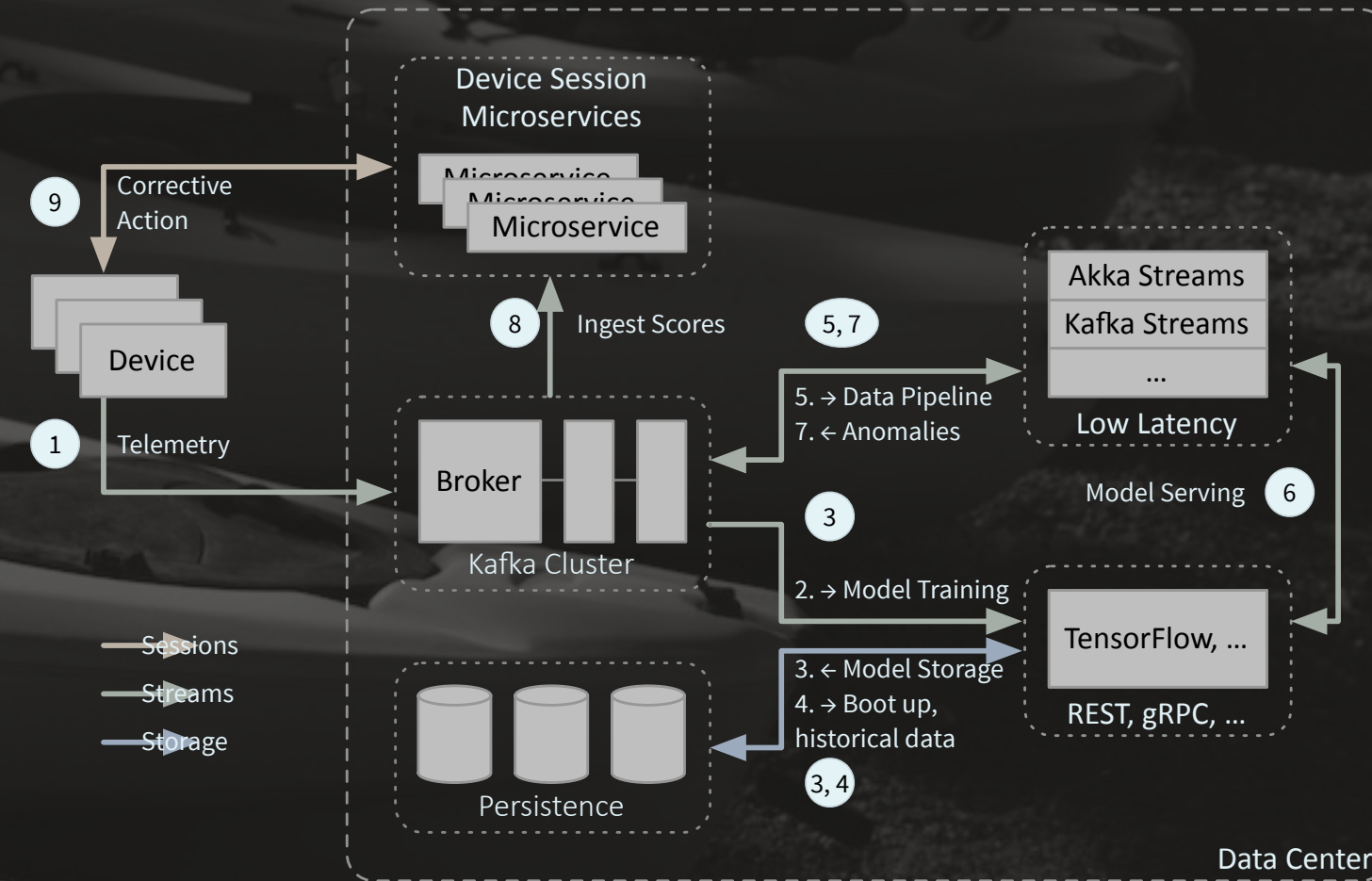


TensorFlow Serving



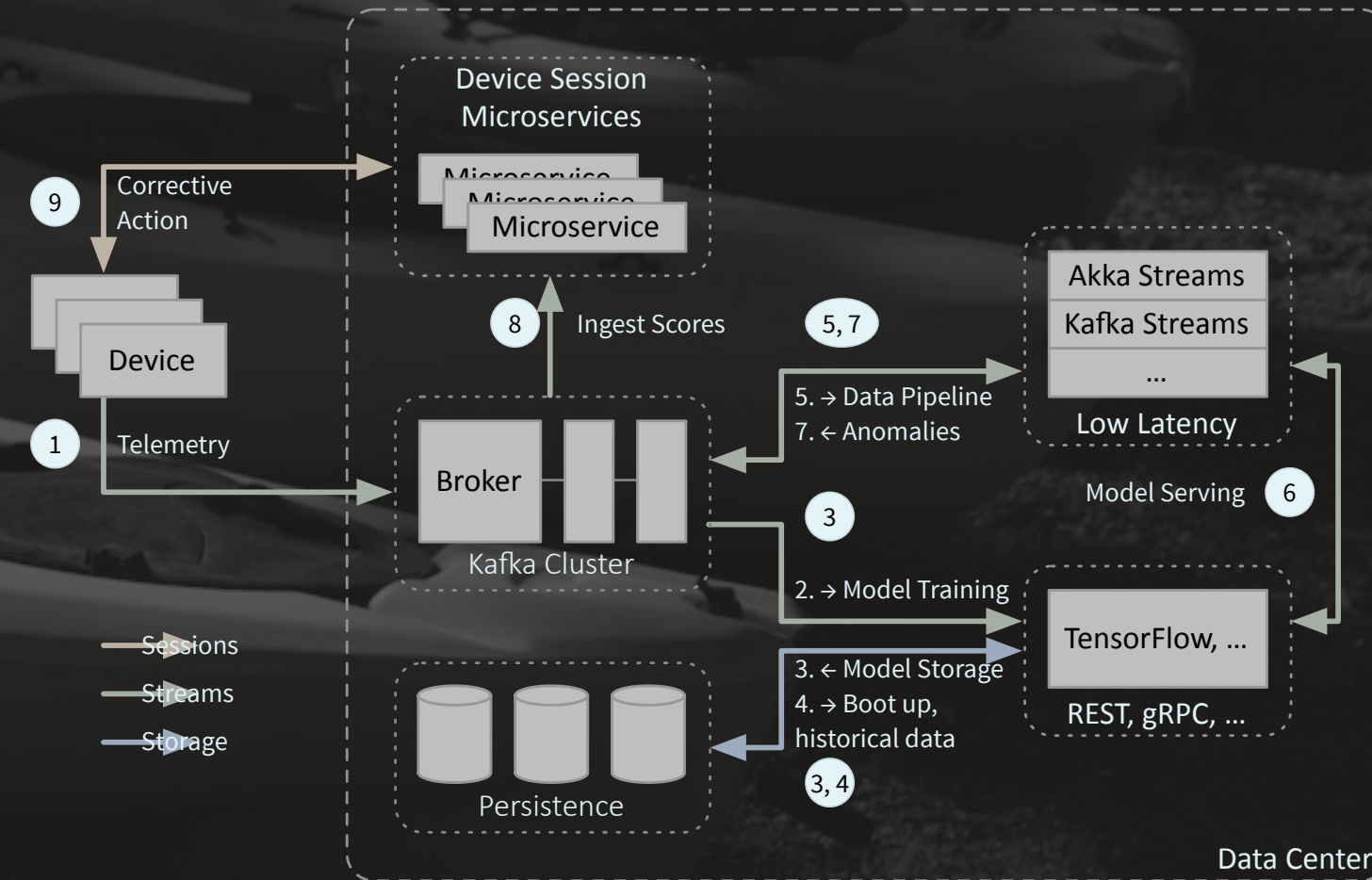
Model Serving as a Service

- Pros:
 - A familiar integration pattern
 - Decouples “concerns”: AI tools, scaling, upgrading, ...
 - One system for training and scoring



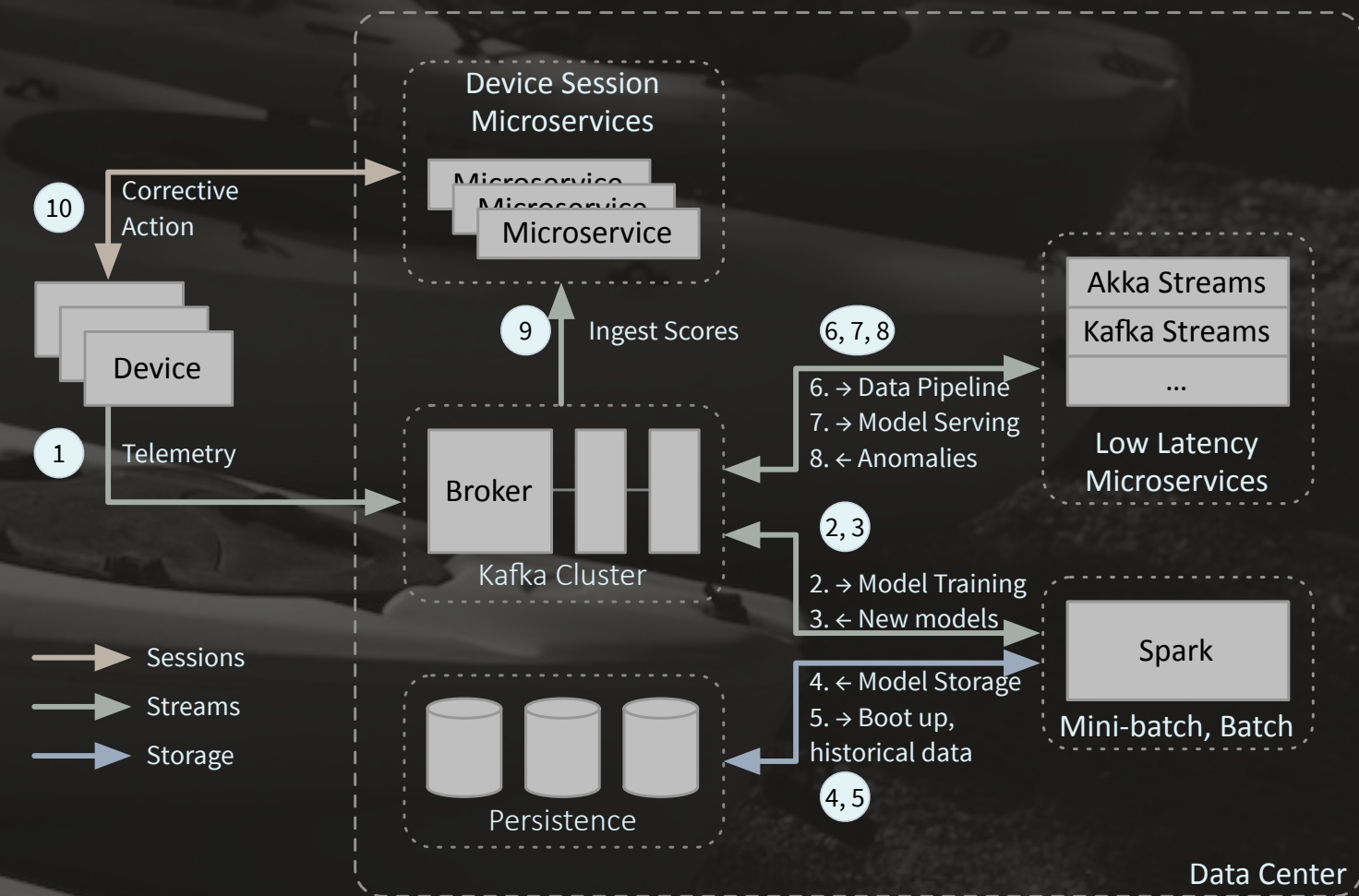
Model Serving as a Service

- Cons:
- Overhead of invocation, e.g., REST
- ML Pipeline becomes a unique production workflow



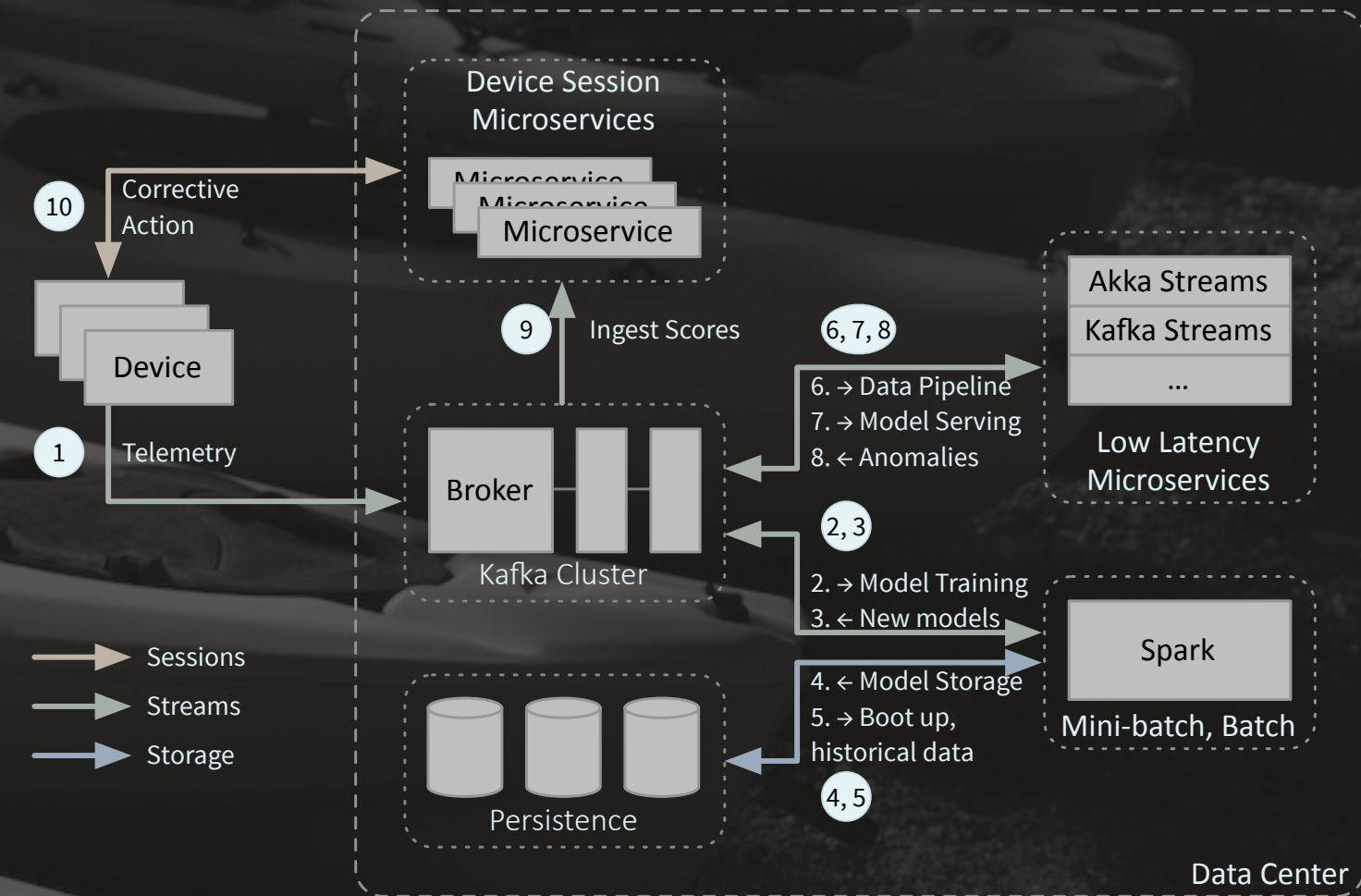
Embedded Model Serving

- Pros:
 - Lowest scoring overhead - interprocess communication only used for model updates
 - Performance tuning focuses on one system, the data pipeline



Embedded Model Serving

- Cons:
- Model parameters must be serialized
- More complexity
- Model serving library must be “compatible” with training system





Updating Models in Production

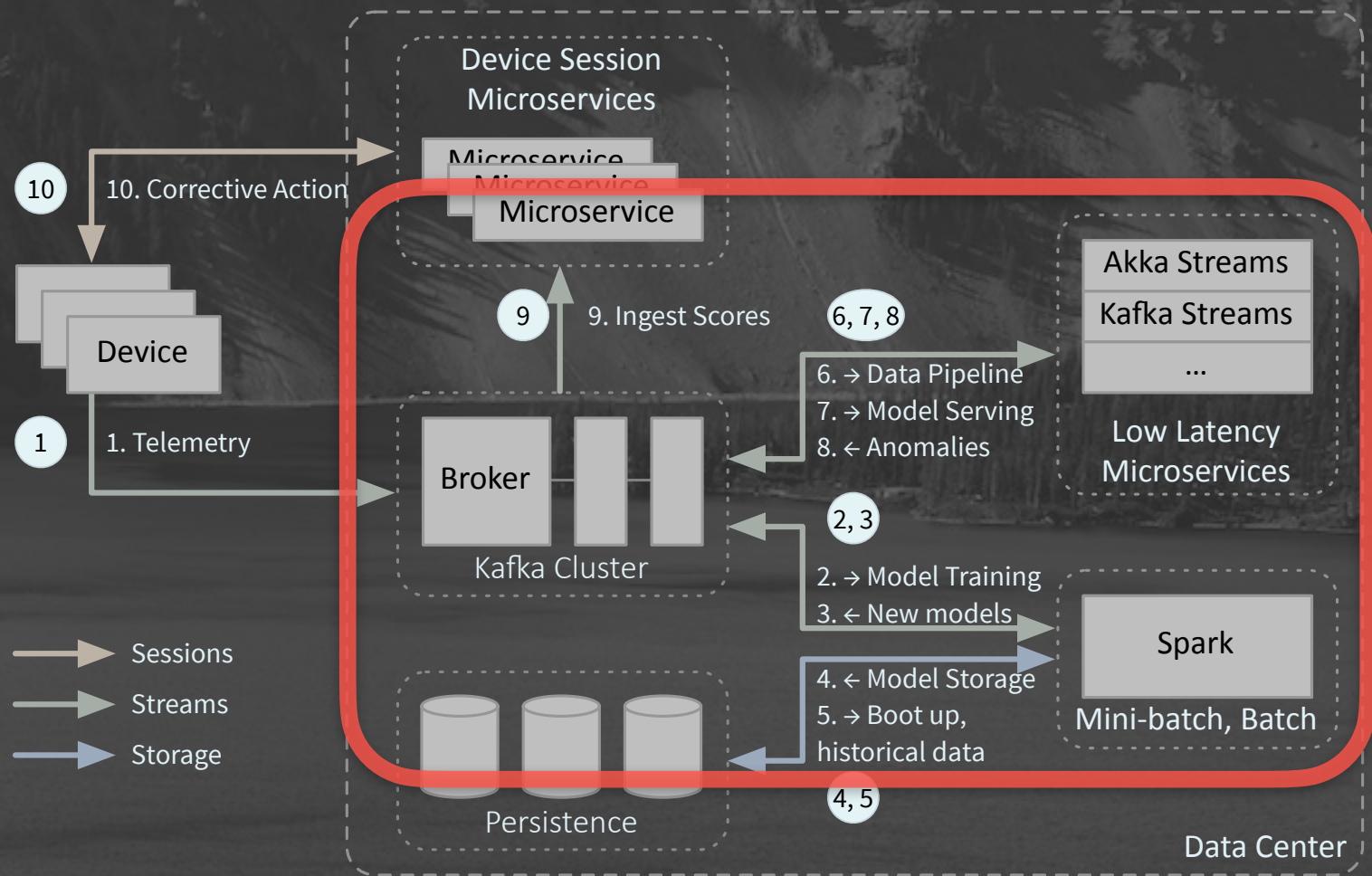
@deanwampler

Model Updates

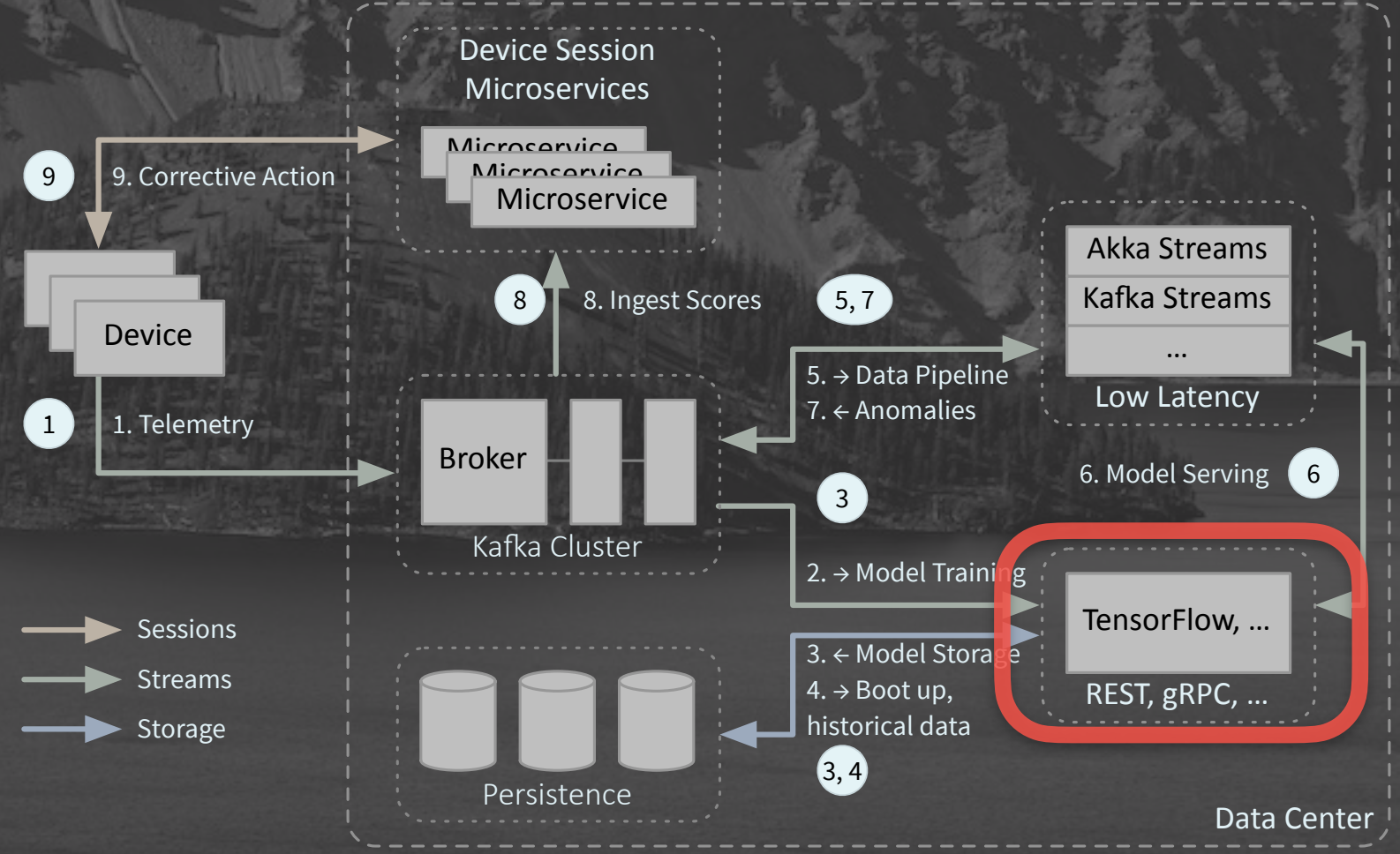
- Concept Drift - models grow stale
 - They have a half life, too
- So, periodically retrain, then serve the new model, ideally without downtime

Retraining Considerations

- How do you measure model quality?
- What's the trade-off between model performance vs. retraining cost?
- How far back in the data set do you go when training?



Complex to update
embedded models!



Model updates can be
straightforward

Auditing

- Kind of Model
- Parameters and hyperparameters
- When trained
- Data used for training
- When deployed, undeployed, etc.
- ...

Auditing

- Quality metrics
- Serving metrics (how many records, scoring times...)
- Provenance of decision to retrain
 - The metrics gathered above that were used to decide when to retrain



Dusty Milky Way

Mars last Summer

References

- Ideas:
 - [Ben Lorica on 9 AI Trends](#)
 - [Paco Nathan's Data Governance Talk](#)

You can get these slides with the links here:
polyglotprogramming.com/talks

References

- Ideas:
 - O'Reilly Radar: [Data](#), [AI](#), [others](#)
 - [distill.pub](#)
 - [The Algorithm](#)
 - [The Gradient](#)

References

- A few research papers, etc.
 - Incremental training
 - an example
 - Continual learning
 - Explainability

References

- [Kubeflow](#)
- [MLFlow](#)
- [DVC](#)
- [AWS SageMaker](#)
- [Fiddler](#) (explainable AI)

References

- General Information about Stream Processing
 - [My O'Reilly Report on Architectures](#)
 - [Streaming Systems Book](#)
 - [Stream Processing with Apache Spark](#)
 - [Designing Data-Intensive Apps book](#)

References

- Other Talks
 - [Strata Talk on ML in a Streaming Context](#)
 - [Stream All the Things! \(video\)](#)
 - [Streaming Microservices with Akka Streams and Kafka Streams \(video\)](#)

References

- Tutorials
 - [Model serving in streams](#)
 - [Stream processing with Kafka and microservices](#)

Questions?

@deanwampler

dean@deanwampler.com

polyglotprogramming.com/talks

