

SMACK STACK AND BEYOND

BUILDING FAST DATA PIPELINES

Jörg Schad, Mesosphere

@dcos @joerg_schad



Jörg Schad

Software Engineer @Mesosphere



@joerg_schad



@joerg.mesosphere

Ancient Times...



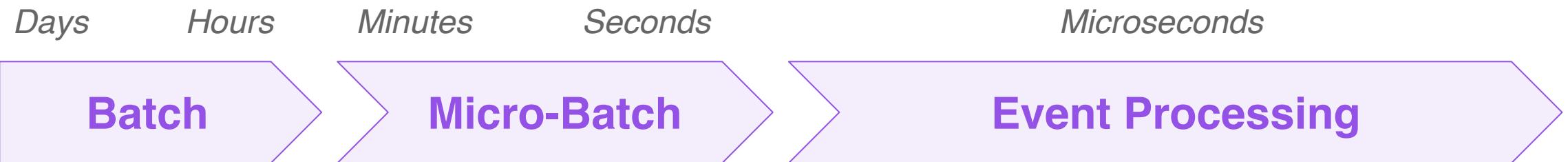
MapReduce is
crunching Data

Today...



We need to turn faster!

Data Processing



Reports what has happened using descriptive analytics

Solves problems using predictive and prescriptive analytics

Billing, Chargeback



Product recommendations



Real-time Pricing and Routing



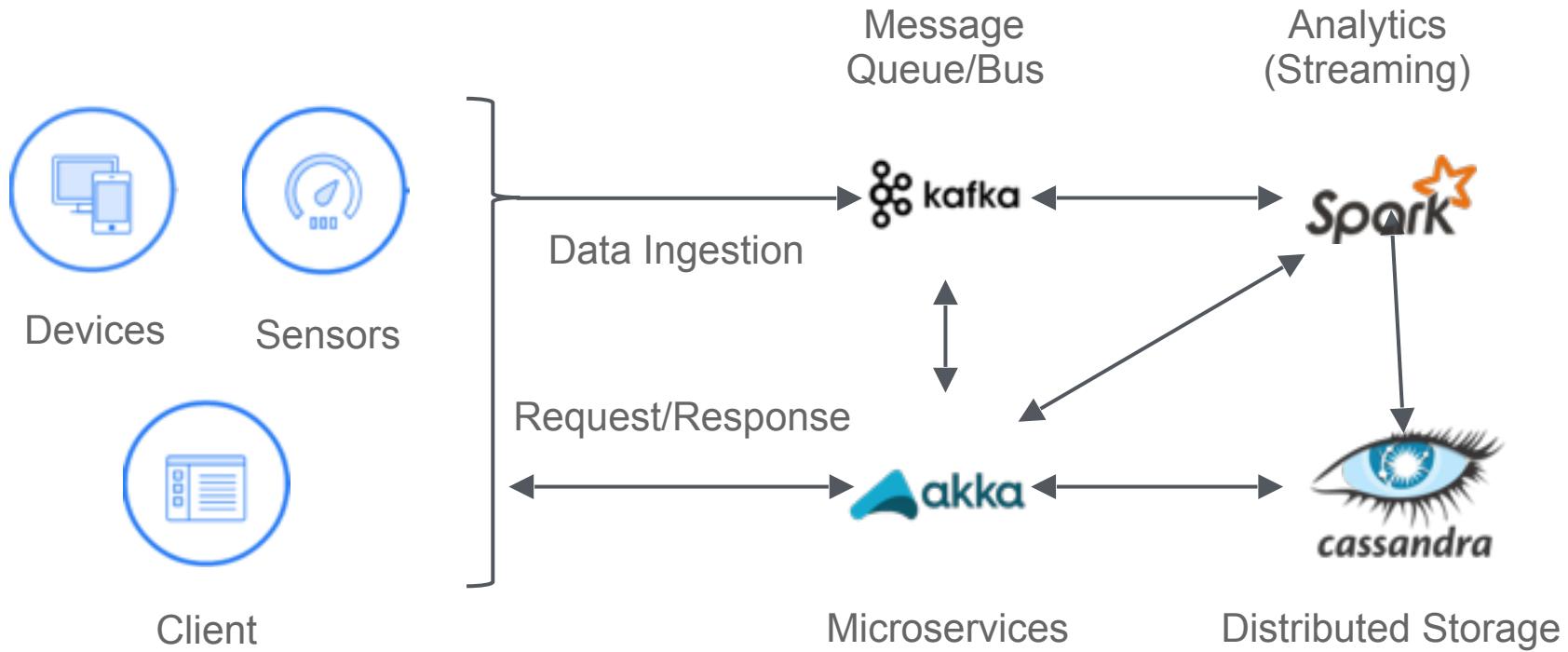
Real-time Advertising



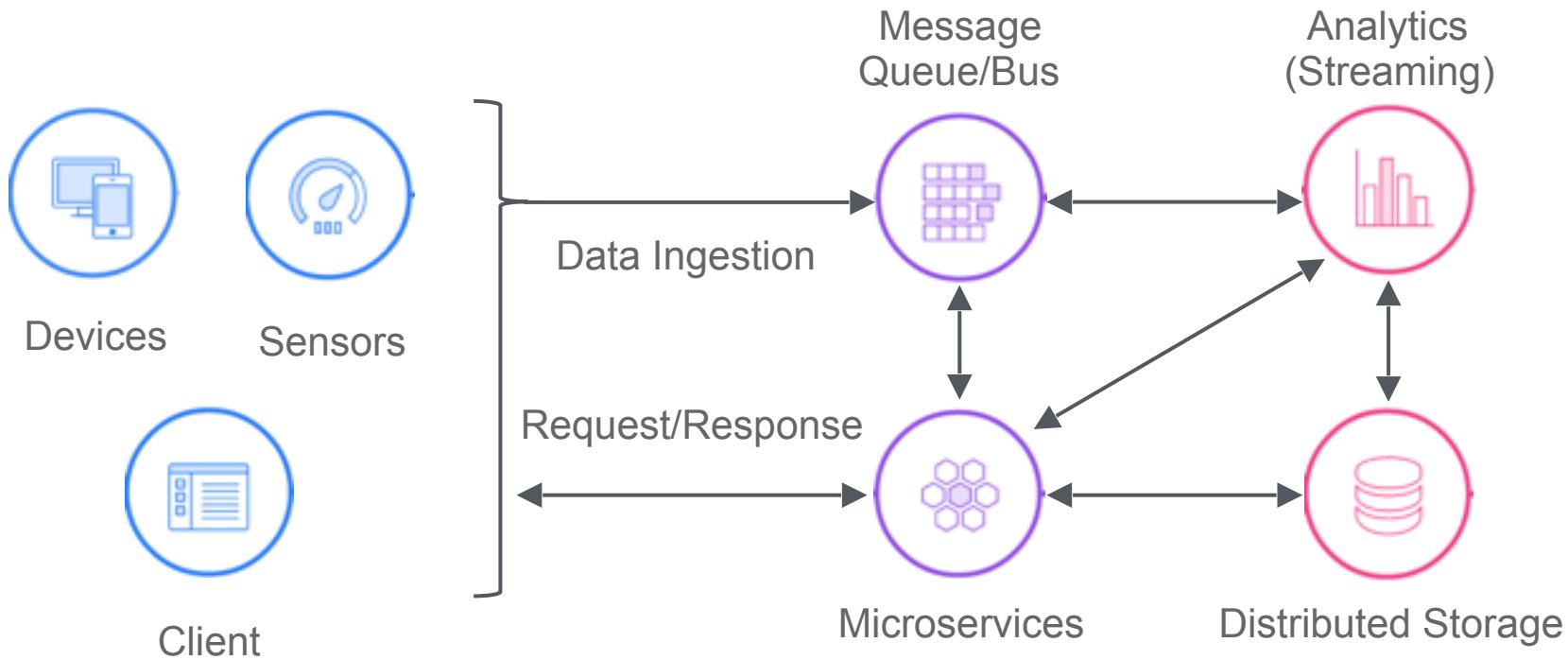
Predictive User Interface



Fast Data Pipelines



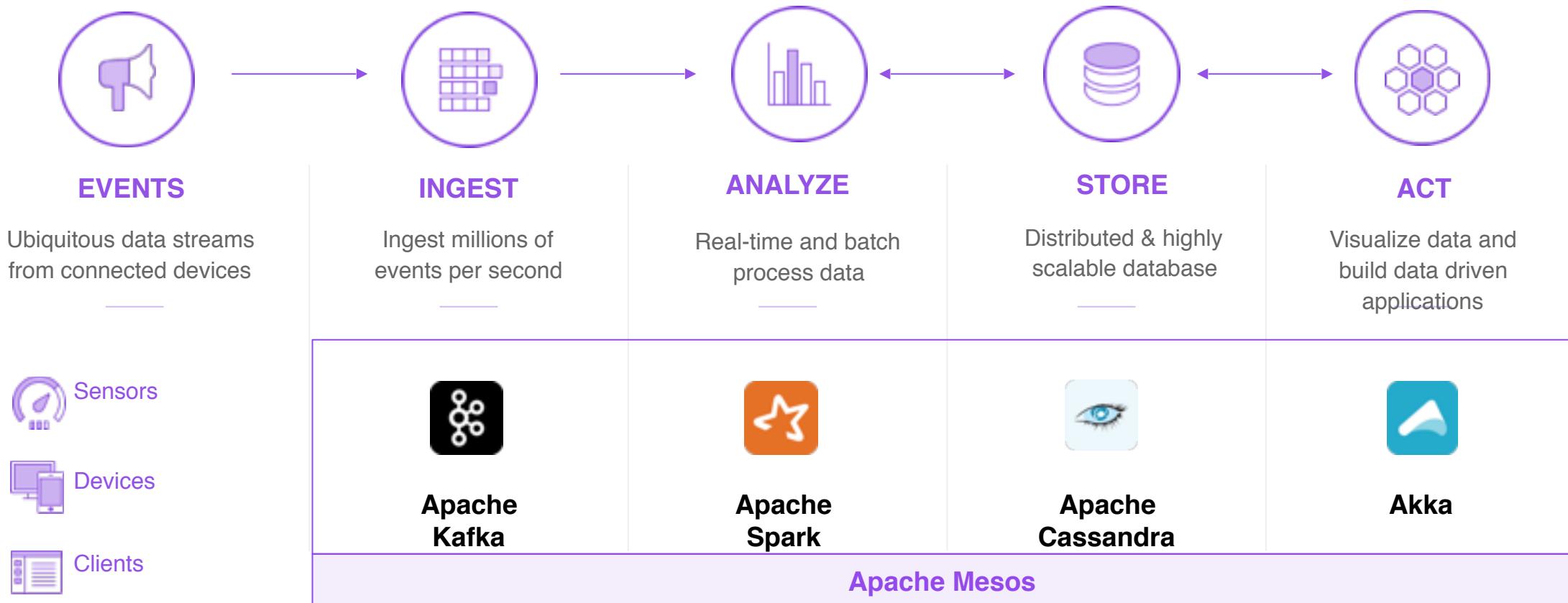
Fast Data Pipelines



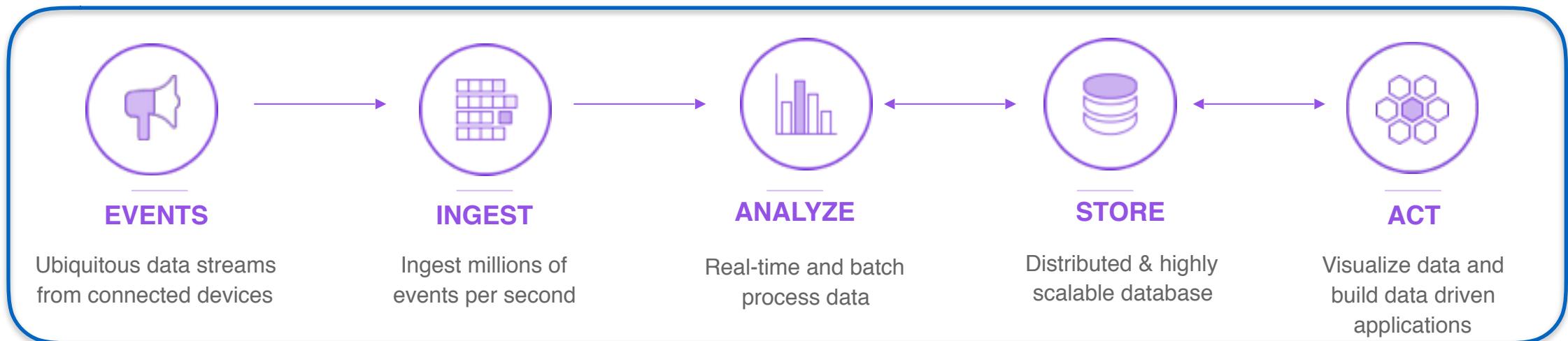
Fast Data Pipelines



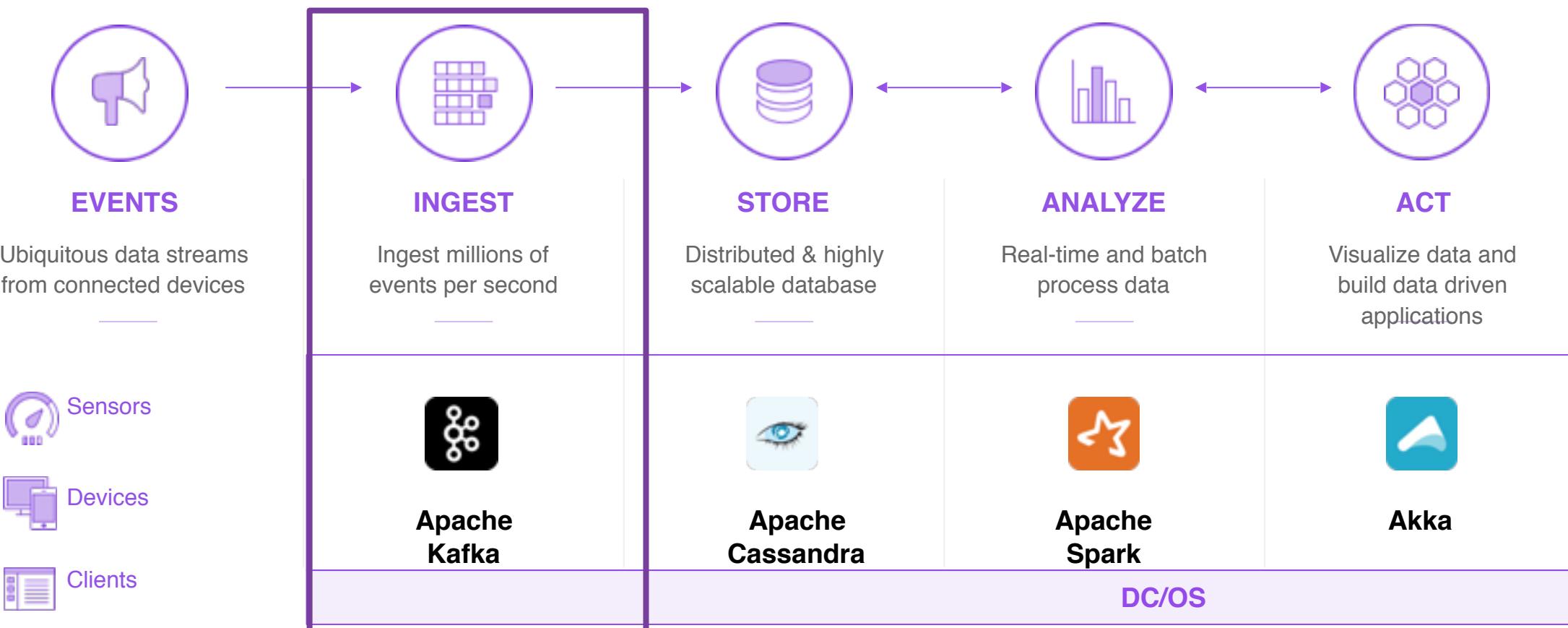
SMACK Stack



SMACK Stack

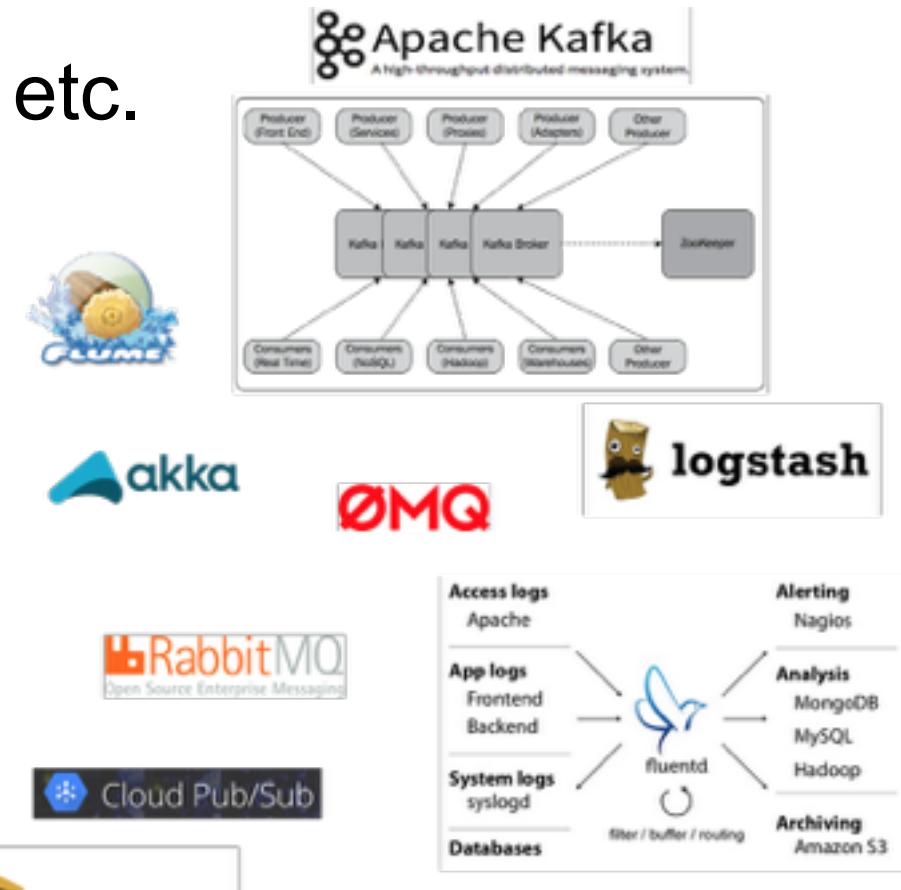


Message Queues



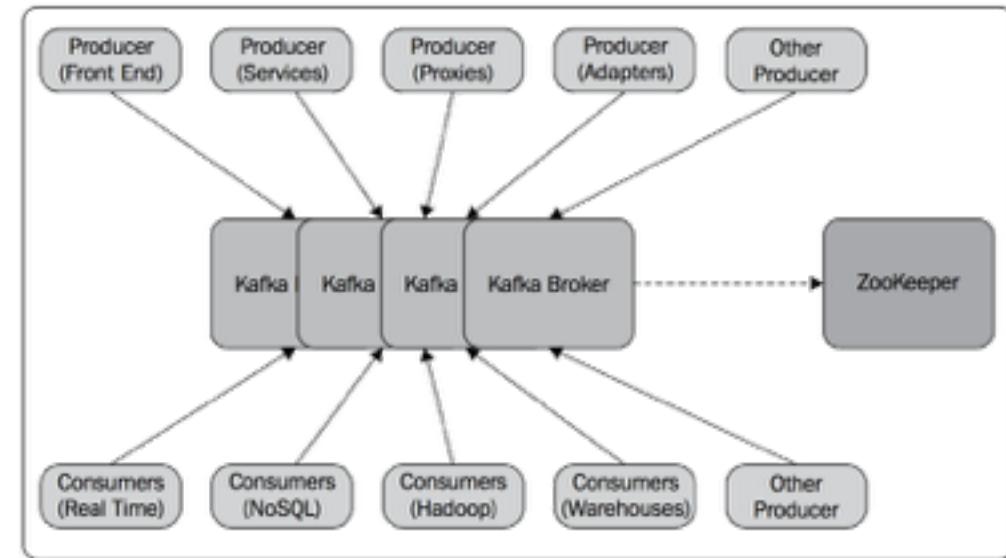
MESSAGE QUEUES

- Apache Kafka
- ØMQ, RabbitMQ, Disque (Redis-based), etc.
- fluentd, Logstash, Flume
- Akka streams
- cloud-only:
 - AWS SQS
 - Google Cloud Pub/Sub

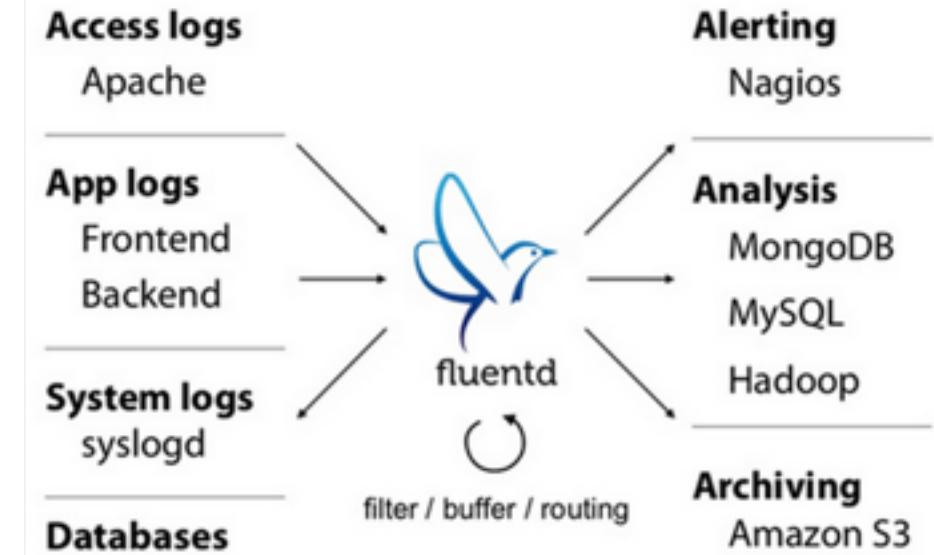
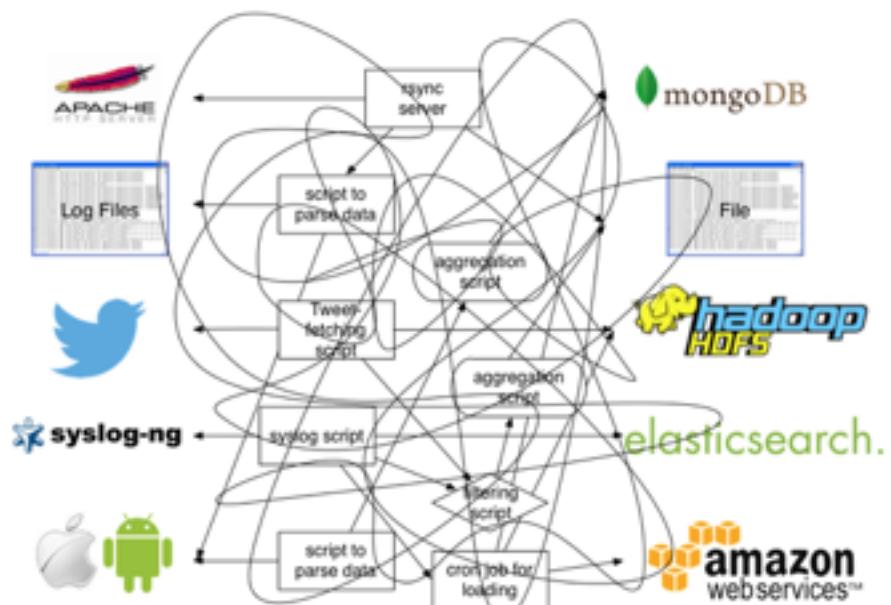


APACHE KAFKA

- High-throughput, distributed, persistent publish-subscribe messaging system
- Originates from LinkedIn
- Typically used as buffer/de-coupling layer in online stream processing



fluentd



HOW TO CHOOSE?

- Scalability
- Message Type
 - Log vs ...
- Delivery Guarantees/Message durability
- Routing Capabilities
- Failover
- Community
- Mesos Support ;-)

DELIVERY GUARANTEES

At most once—Messages may be lost but are never redelivered.

At least once—Messages are never lost but may be redelivered.

Exactly once—this is what people actually want, each message is delivered once and only once.

Murphy's Law of Distributed Systems:

Anything that can go wrong, will go wrong ... partially!

Routing

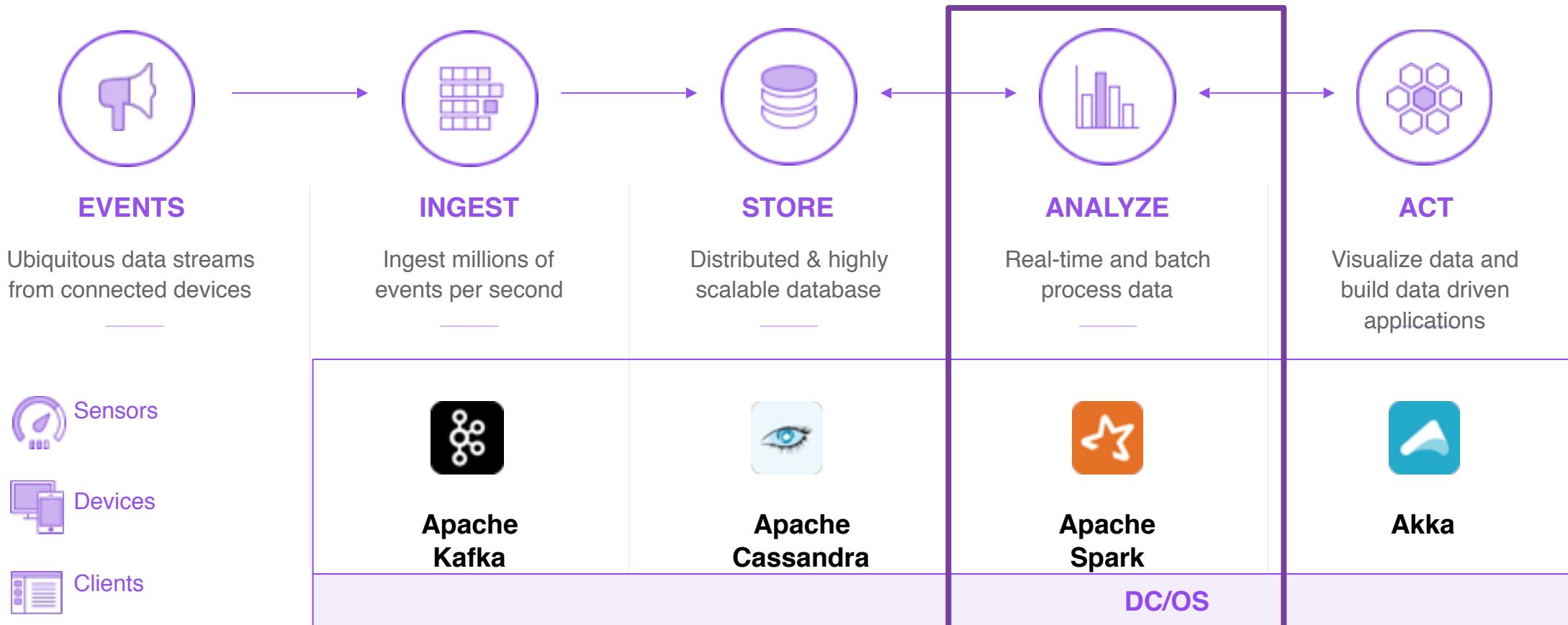
Simple Pipes



Routing



Stream Processing



STREAM PROCESSING

- Apache Storm
- Apache Spark
- Apache Samza
- Apache Flink
- Apache Apex
- Concord
- cloud-only: AWS Kinesis,
Google Cloud Dataflow



APACHE SPARK

Spark SQL

Spark Streaming

MLlib
(machine learning)

GraphX
(graph processing)

Spark core (RDD)

Mesos

Standalone

YARN

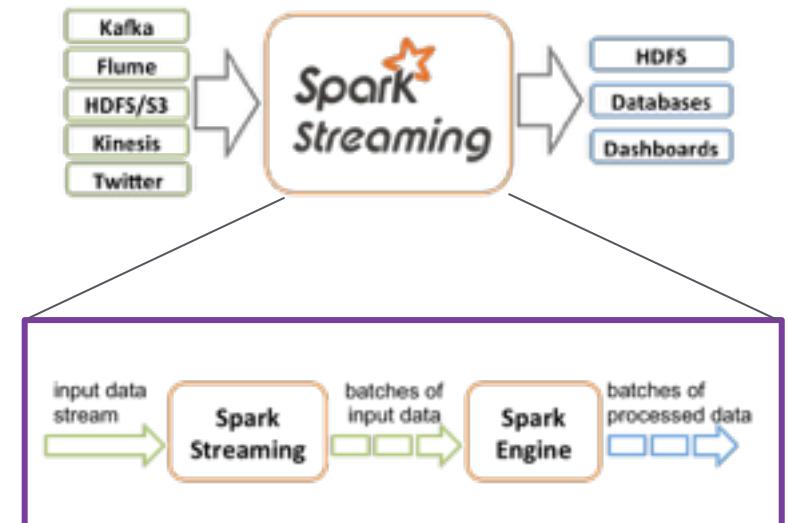
Filesystem (local, HDFS, S3) or data store (HBase, Cassandra, Elasticsearch, etc.)

APACHE SPARK (STREAMING 2.0)

Typical Use: distributed, large-scale data processing;
micro-batching

Why Spark Streaming?

- Micro-batching creates very low latency, which can be faster
- Well defined role means it fits in well with other pieces of the pipeline

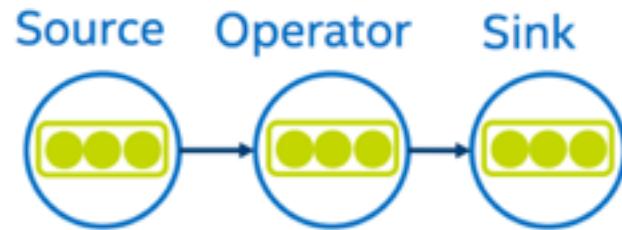


HOW TO CHOOSE?

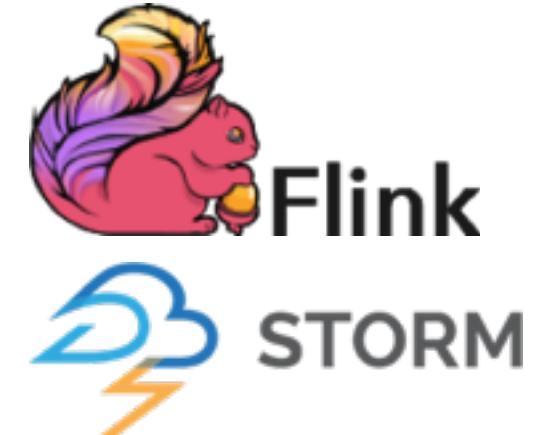
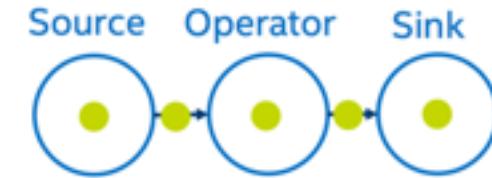
- Execution Model
 - Native Streaming vs Microbatch
- Fault Tolerance Granularity
 - Per record, per batch
- Delivery Guarantees
- API
 - SQL
 - Spark
- Performance....
 - Realtime ≠ Realtime
- Community
- Mesos Support ;-)

EXECUTION MODEL

Micro-Batching

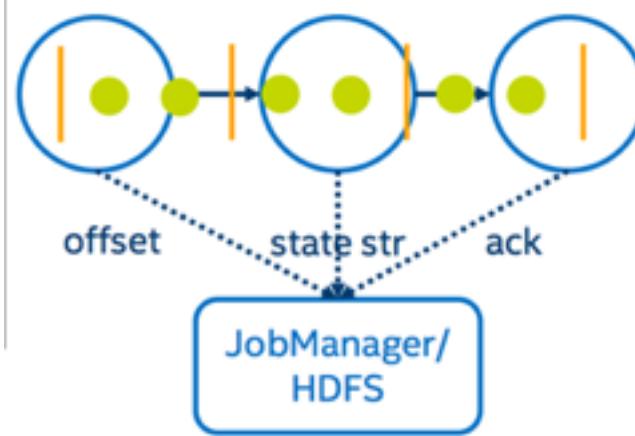
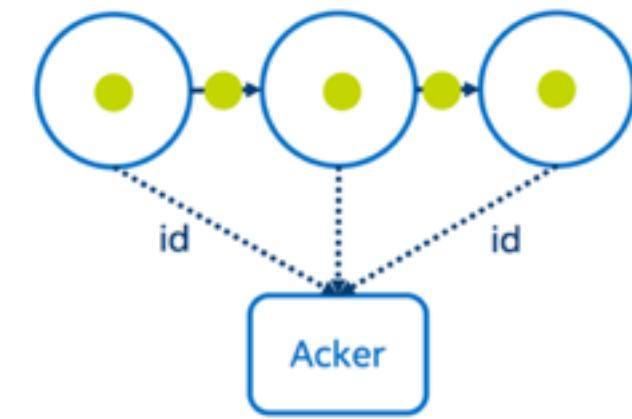


Native Streaming

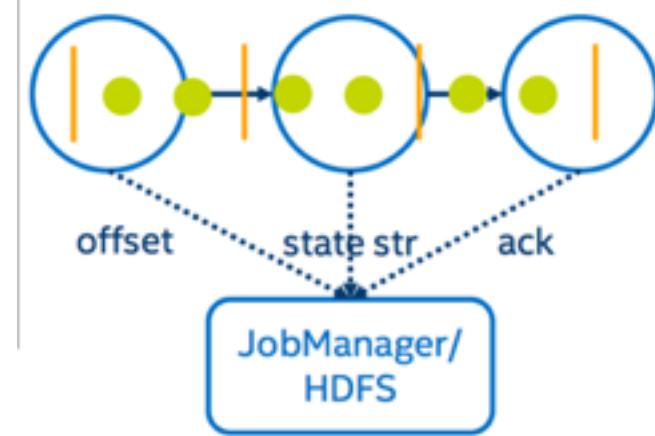


FAULT TOLERANCE

Checkpoint per “Batch”
Ack-Per-Record



Checkpoint per Batch

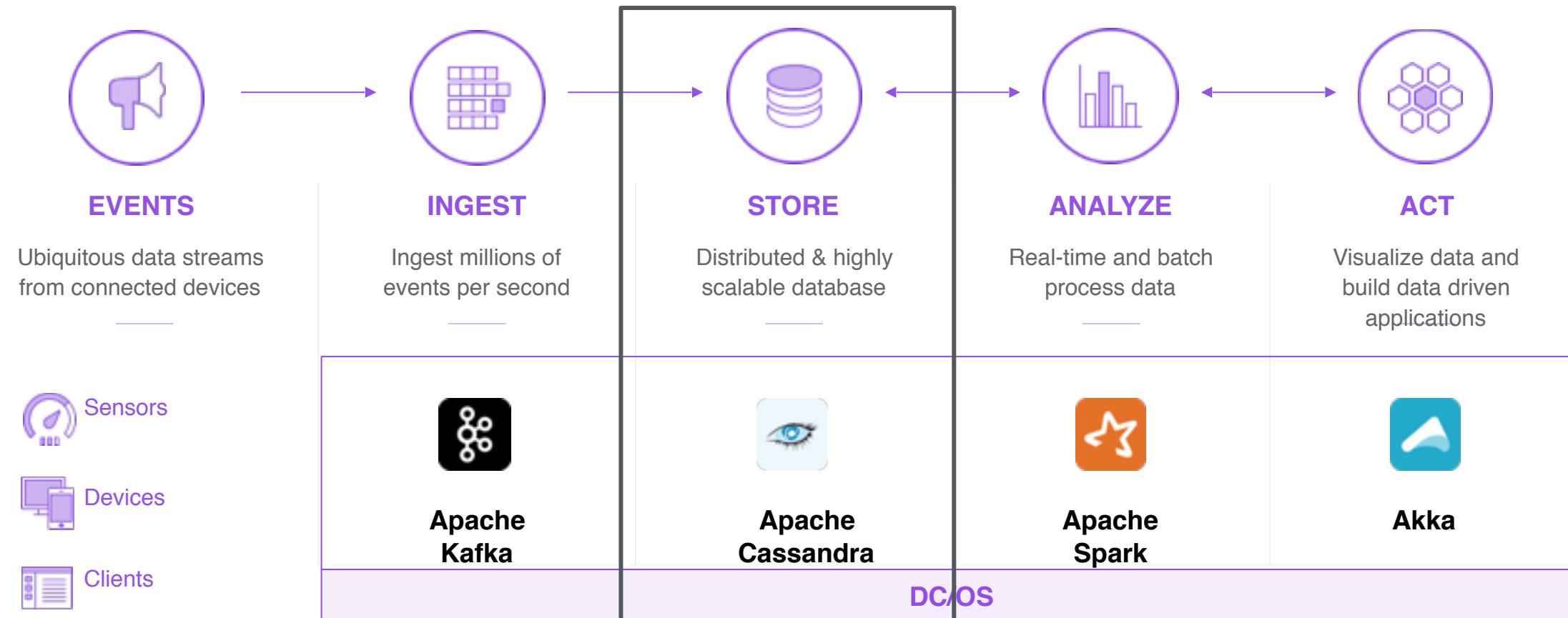


DELIVERY GUARANTEES

“Exactly once”
At least Once



Storage



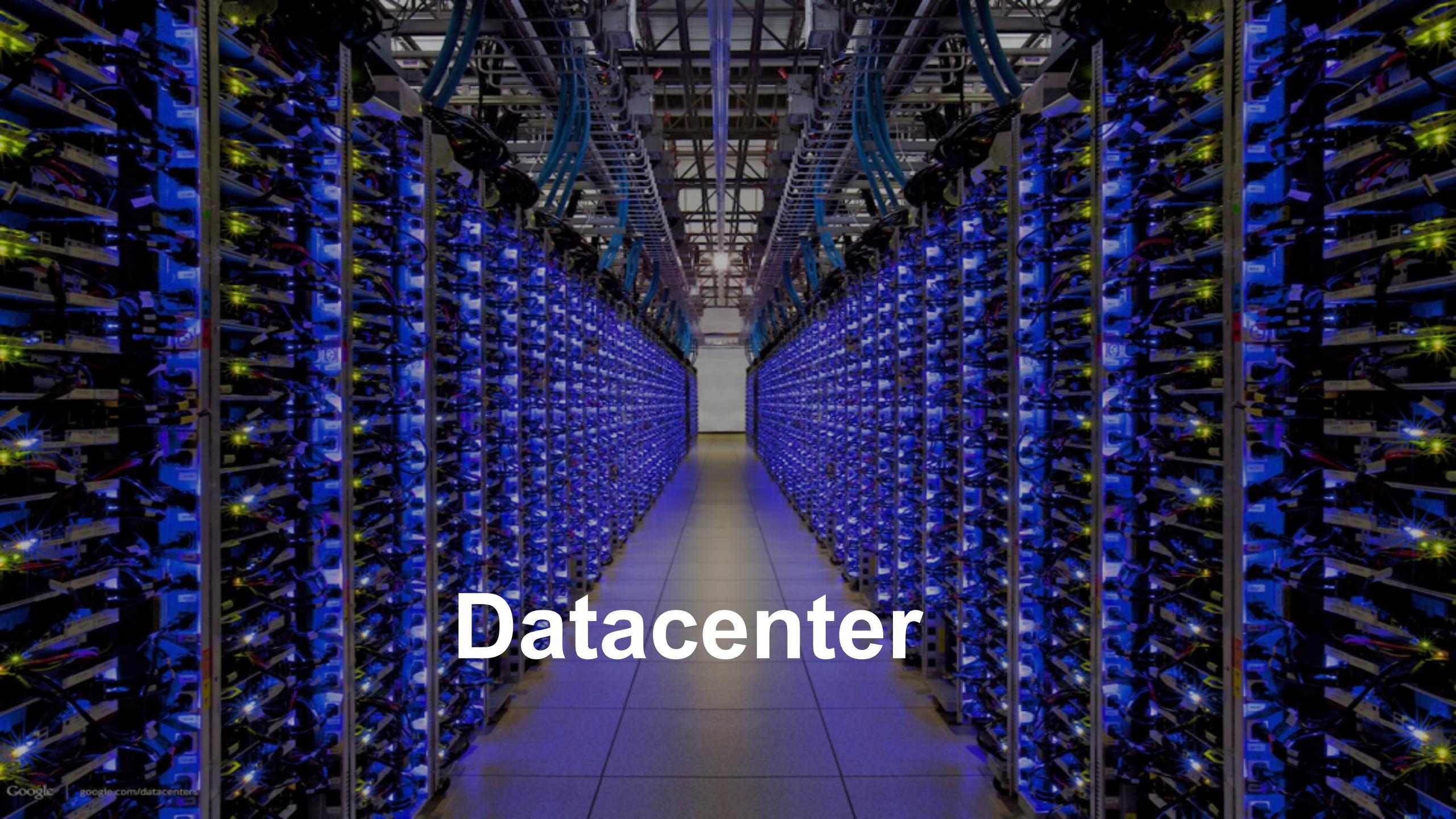
Datastores



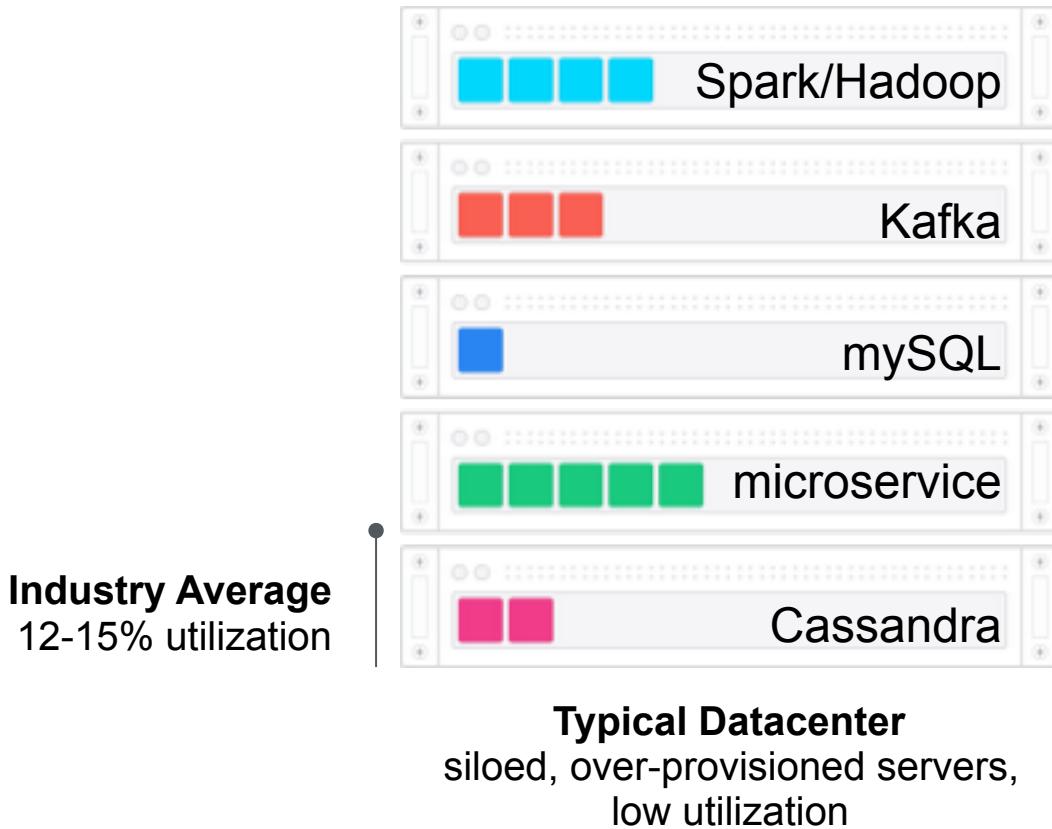
Data Model

Relational	Key-Value	Graph	Document	Time-Series	Files
<ul style="list-style-type: none">• Schema• SQL• Foreign Keys/Joins• OLTP/OLAP   	<ul style="list-style-type: none">• Simple• Scalable• Cache  	<ul style="list-style-type: none">• Complex relations• Social Graph• Recommendation• Fraud detections  	<ul style="list-style-type: none">• Schema-Less• Semi-structured queries• Product catalogue• Session data  		  

Datacenter

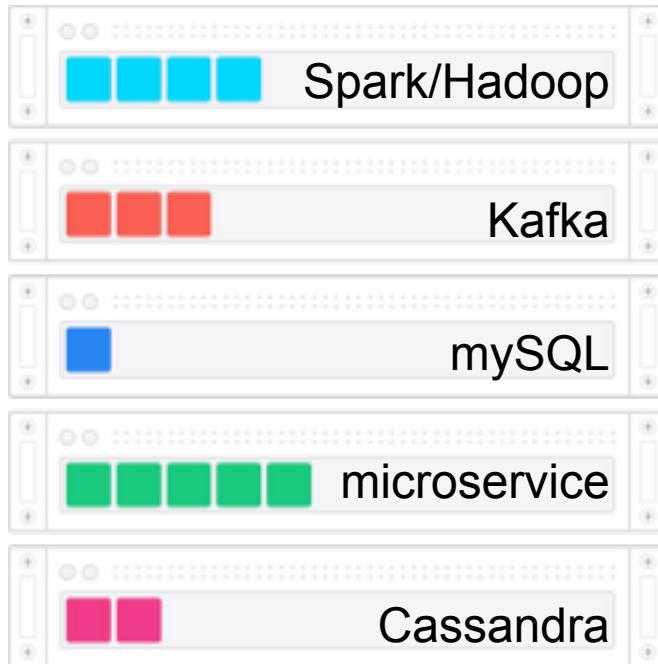


NAIVE APPROACH





MULTIPLEXING OF DATA, SERVICES, USERS, ENVIRONMENTS



Typical Datacenter
siloed, over-provisioned servers,
low utilization



Mesos/ DC/OS
automated schedulers, workload multiplexing onto the
same machines

Apache Mesos

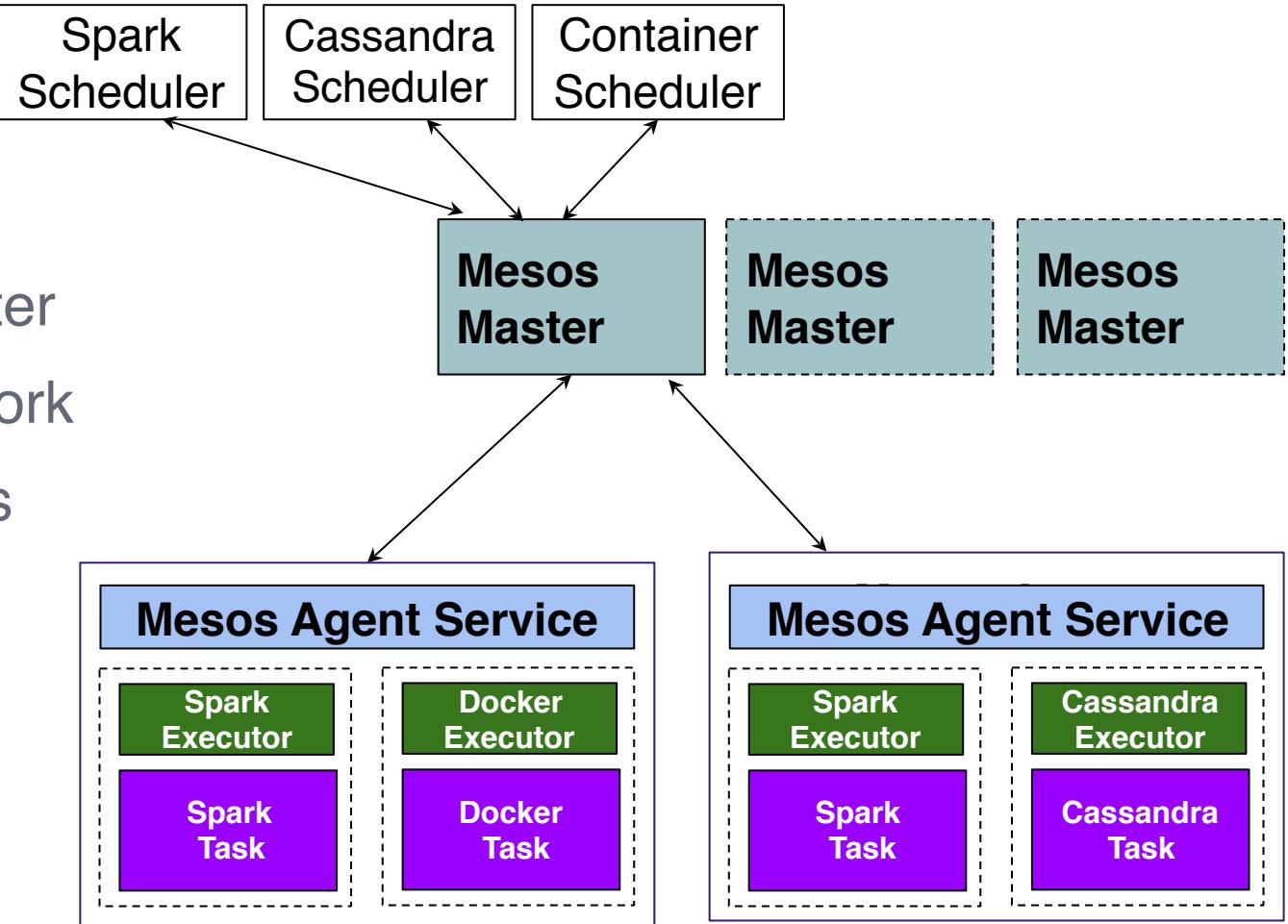
- A top-level Apache project
- A cluster resource negotiator
- Scalable to 10,000s of nodes
- Fault-tolerant, battle-tested
- An SDK for distributed apps
- Native Docker support



MESOS: FUNDAMENTAL ARCHITECTURE

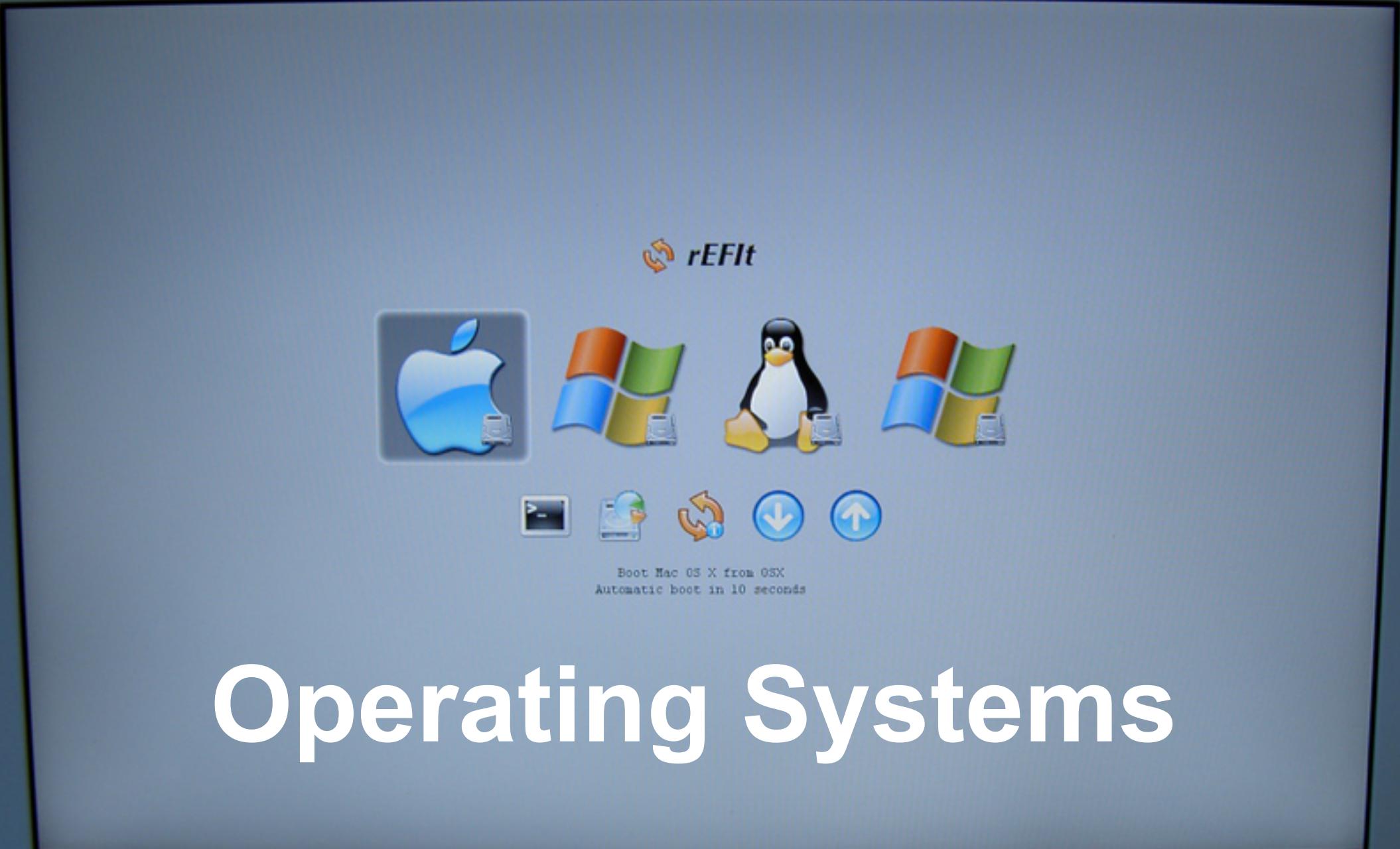
Two-level Scheduling

1. Agents advertise resources to Master
2. Master offers resources to Framework
3. Framework rejects / uses resources
4. Agent reports task status to Master



Challenges

- Mesos is just the kernel
- Need for OS:
 - Scheduler
 - Monitoring
 - Security
 - CLI
 - Package Repository
 - ...

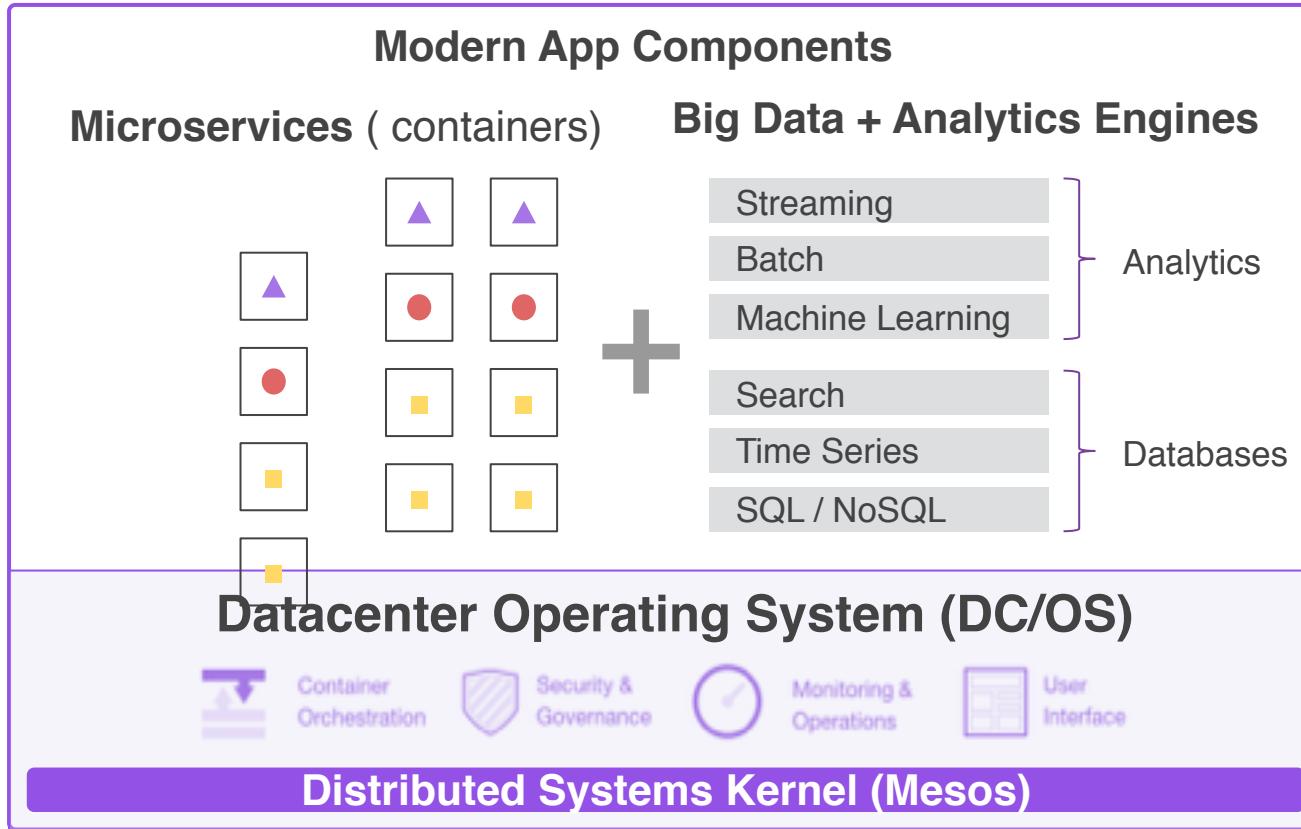


Operating Systems



DC/OS

DC/OS



DC/OS

- Container operations & big data operations
- Security, fault tolerance & high availability
- Open Source (ASL2.0)
- Based on Apache Mesos
- Production proven at scale

DC/OS Universe

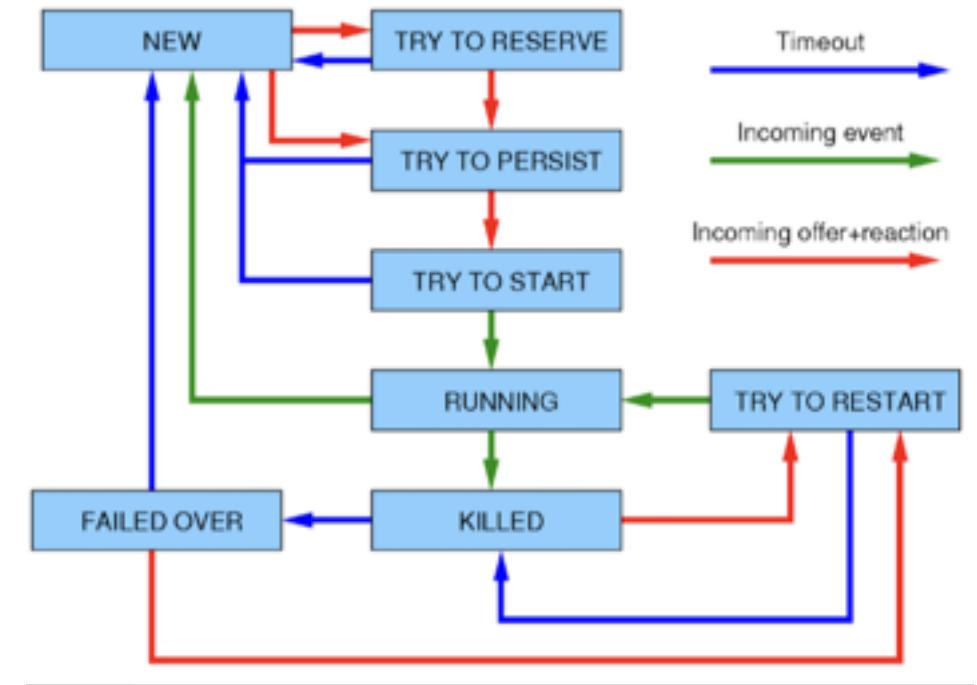
- Datacenter-wide services to power your apps
- Turnkey installation and lifecycle management

Any Infrastructure

- Requires only a modern linux distro (windows coming soon)
- Hybrid Datacenter

Developing Distributed Services

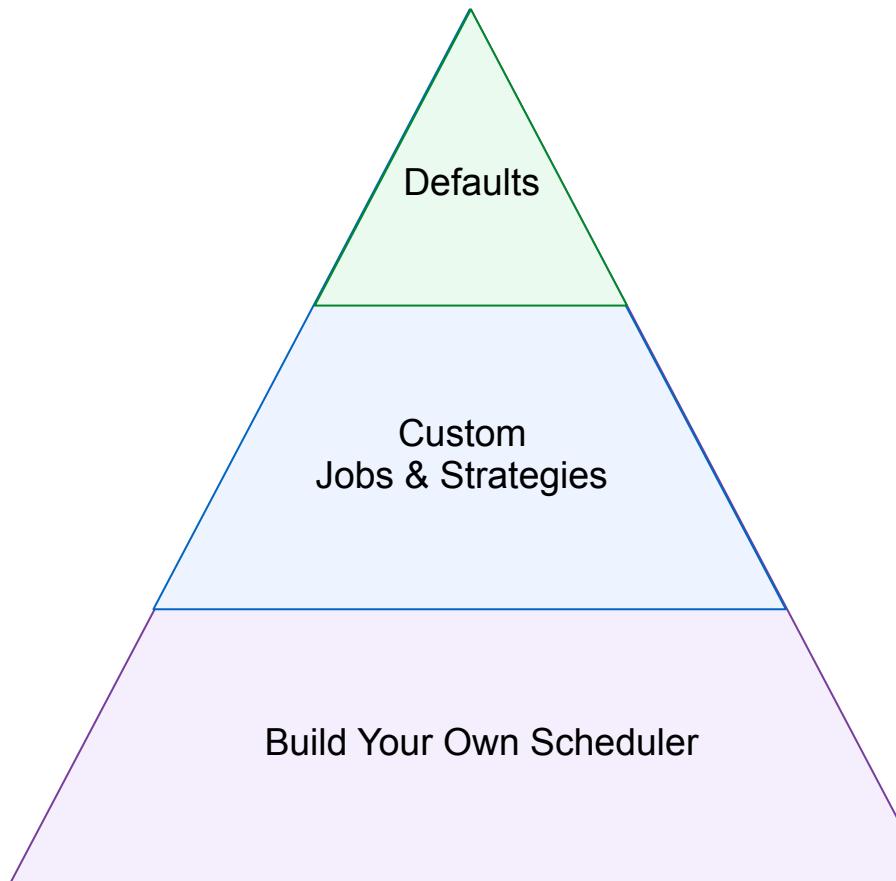
- Failures (Task, Node, Network,...)
- Zero Downtime Upgrades
- Persistence
- Multiple Frameworks
- Service Discovery
- Metrics
-



Operating Distributed Services

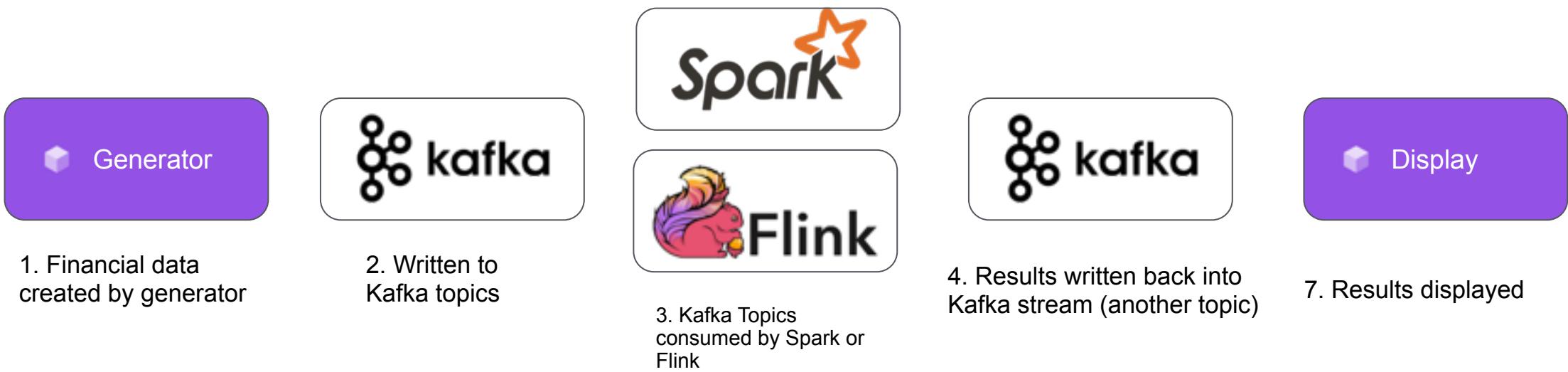


Distributed Services: Challenges



- As simple as Docker Compose
 - Don't need to write any Java code
 - Don't need to be an app expert
-
- Need to be an app expert
 - Need to write a little Java code
 - Don't want to understand DC/OS
-
- Can't use the default scheduler
 - Need to write a lot of Java code
 - Willing to understand DC/OS

Demo Time





Keep it running!

SERVICE OPERATIONS

- Configuration **Updates** (ex: Scaling, re-configuration)
- Binary **Upgrades**
- Cluster **Maintenance** (ex: Backup, Restore, Restart)
- **Monitor** progress of operations
- **Debug** any runtime blockages

S

SPARK
SUMMIT

Questions?

dcos.io
[Demo](#)

@dcos @joerg_schad