Kinesis

Spark

Streaming

# Spark Streaming

# +

# Kinesis

Chris Fregly
AdvancedAWS Meetup
8/12/2014

# Who am I?

Former Netflix'er
(netflix.github.io)



Spark Contributor
(github.com/apache/spark)





Consultant
(fluxcapacitor.com)

Author
(sparkinaction.com)

# Quick Poll

- Spark, Spark Streaming?

- Hadoop, Hive, Pig?

- EMR, Redshift?

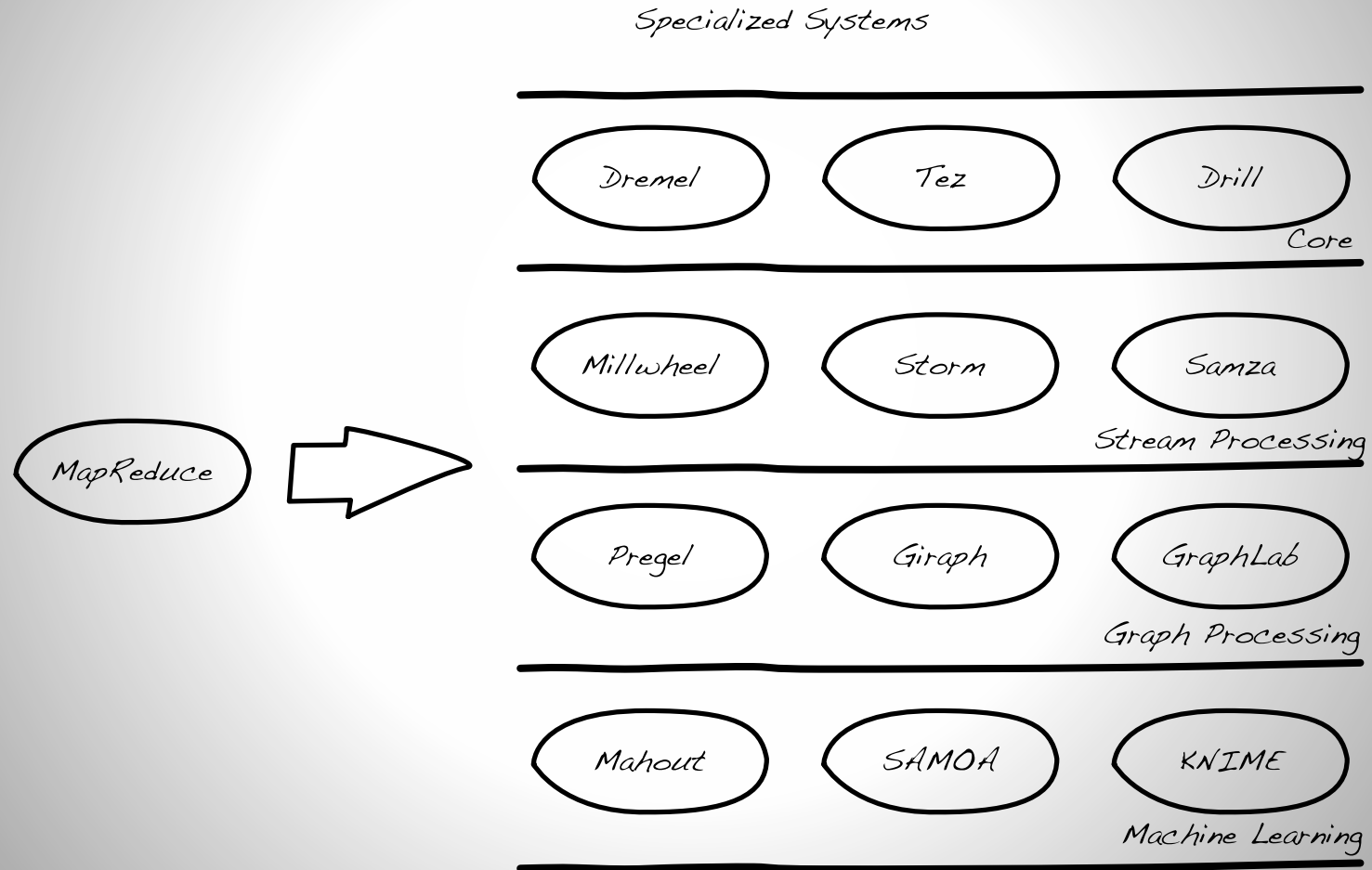- Flume, Kafka, Kinesis, Storm?

- Lambda Architecture?

# Spark Overview

- Part of Berkeley Data Analytics Stack ("badass")
- ~2009, Berkeley AMPLab
- Written in Scala
- Supports Java, Python, SQL, and R
- In-memory whenever possible
- Improved efficiency over MapReduce
  - 100x in-memory, **2-10x on-**disk
- Compatible with Hadoop
  - File formats, SerDes, and UDFs
  - Hive (Shark) and Pig (Spork)

# Spark API

- Richer, more expressive API than just map() and reduce()
  - filter(), join(), distinct(), groupByKey()
- Supports Java, Python, SQL, and R
- Resilient Distributed Dataset (RDD)
  - Core Spark abstraction
  - Partition across cluster
  - Parallel, fault tolerant, immutable, recomputable
- Unified API across all libraries

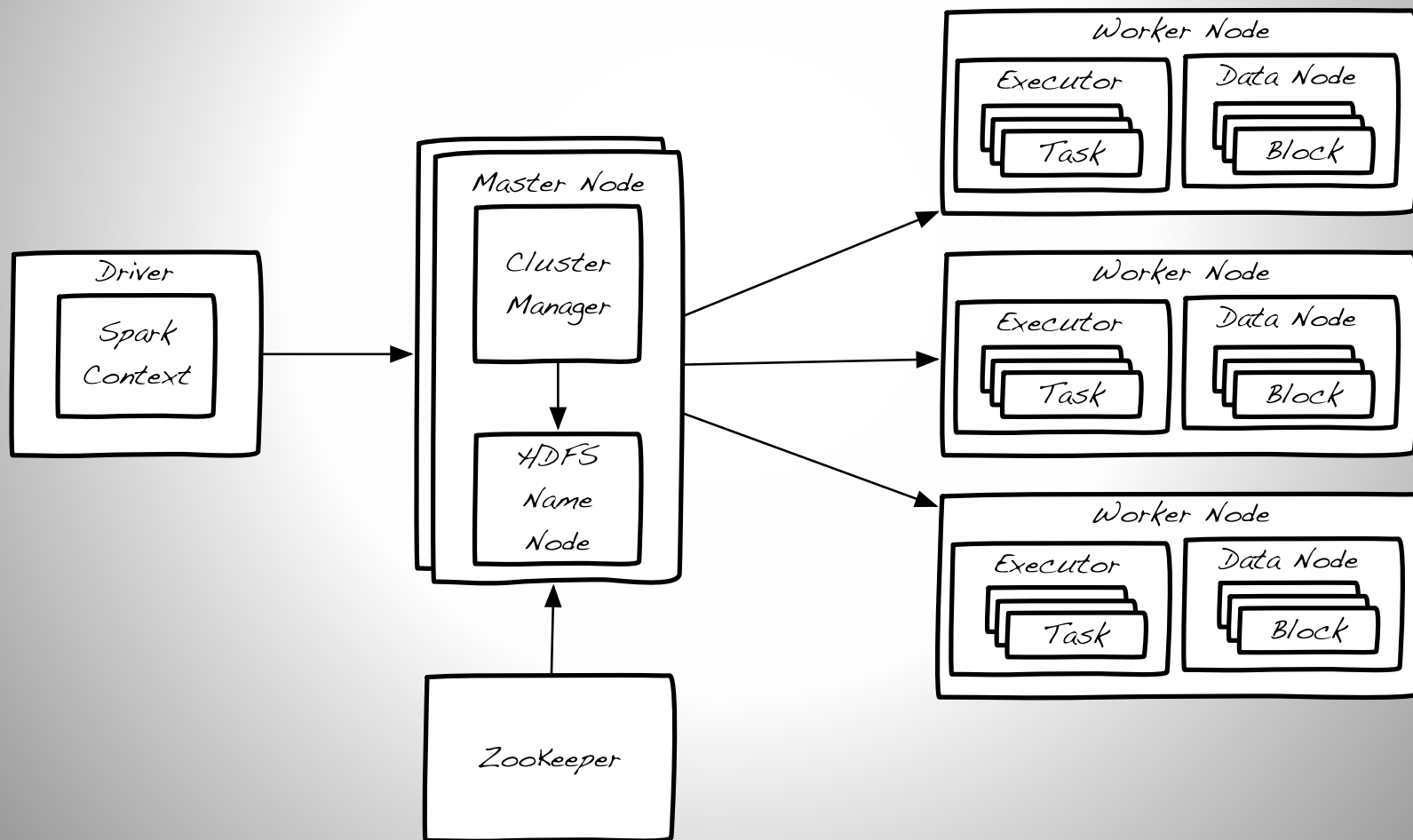# Non-unified Specialized Systems

# Spark Libraries

- Spark SQL (Data Processing)
- Spark Streaming (Streaming)
- MLlib (Machine Learning)
- GraphX (Graph Processing)
- BlinkDB (Approximate Queries)

# Similar Projects

- Spark
  - Microsoft Dryad
  - Apache Tez
  - Impala
  - Google BigQuery
  - Google Cloud Dataflow

# Spark + Hadoop Cluster View



Spark + Hadoop Cluster View

# Master High Availability

- Multiple Master Nodes
- ZooKeeper maintains current Master
- Existing applications and workers will be notified of new Master election
- New applications and workers need to explicitly specify current Master
- Alternatives (Not recommended)
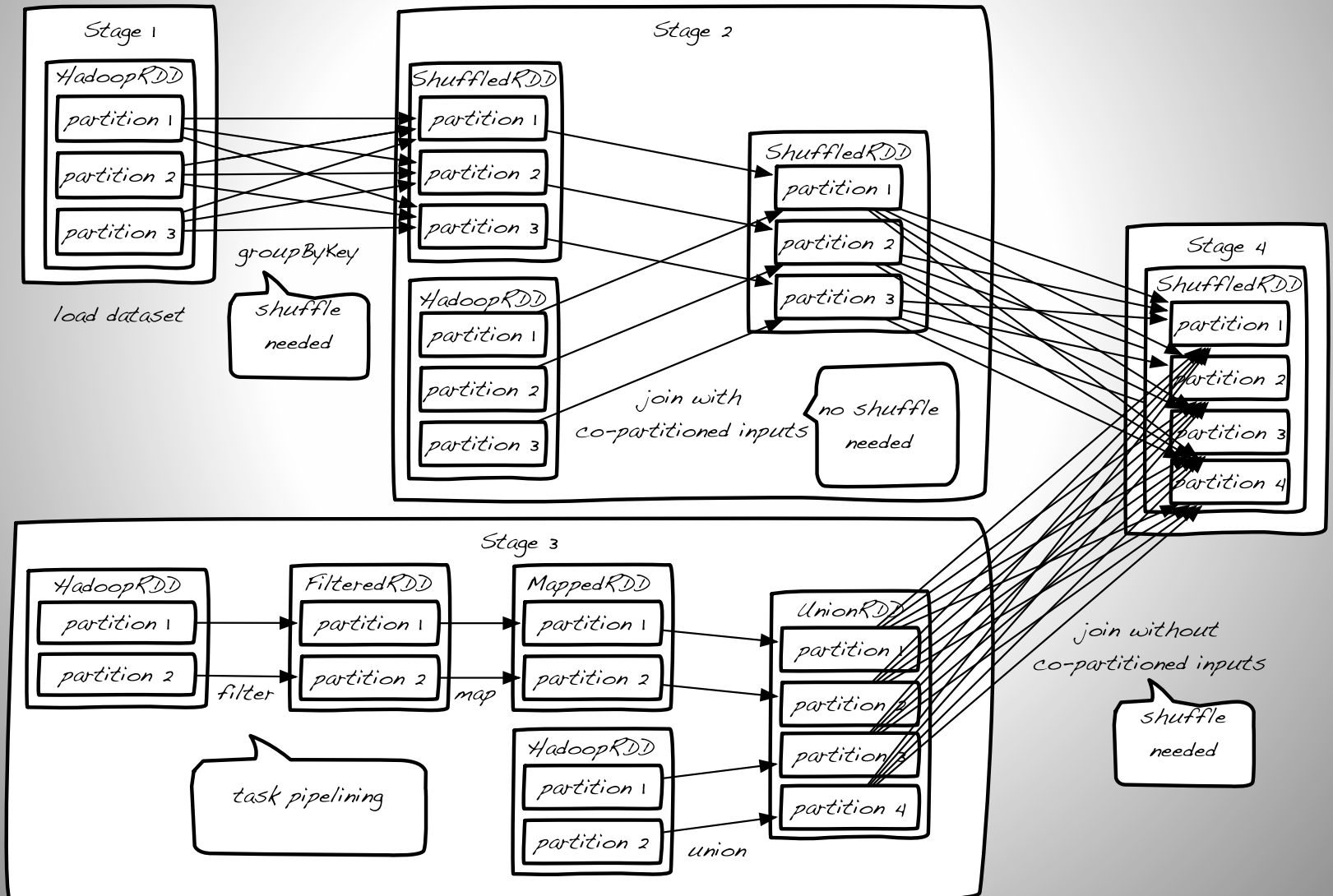  - Local filesystem
  - NFS Mount

# Spark Execution Engine

- General cluster computing engine
- Parallel, distributed, DAG-based
- Lazy evaluation
- Allows optimizations
- Data locality and rack awareness
- Fine-grained fault tolerance using RDD lineage graphs
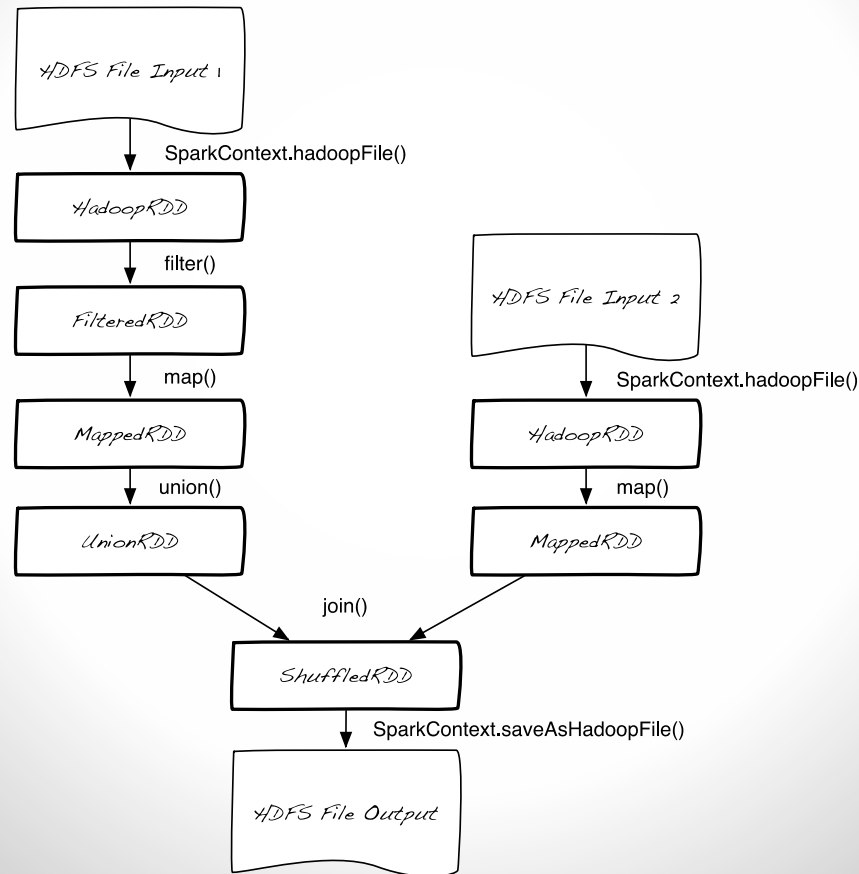
# DAG Scheduler Optimizations

# Lineage-based Fault Tolerance
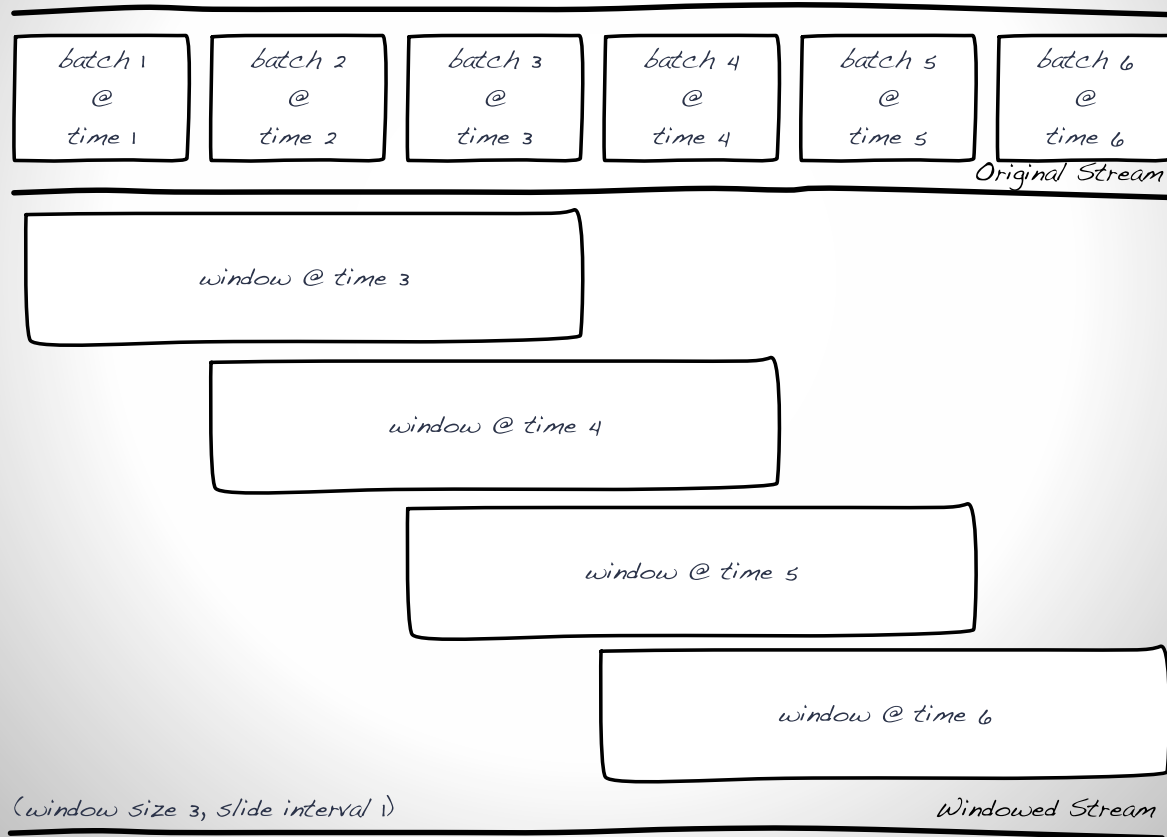
# Spark Streaming Overview

- Low latency, high throughput, fault-tolerant
- DStream:  Micro-batches of RDDs
  - Operations are similar to RDD
  - Lineage for fault-tolerance
- Supports Flume, Kafka, Twitter, Kinesis, etc.
- Built on Spark Core Execution Engine and API
- Long-running Spark Application

# Spark Streaming API

- Rich, expressive API based on core Spark API
    - filter(), join(), distinct(), groupByKey()
- Maintain State
    - updateStateByKey()
- Window Operations
    - Window size & slide interval
- Checkpointing
- Register DStream as a SQL table

# Window Operations
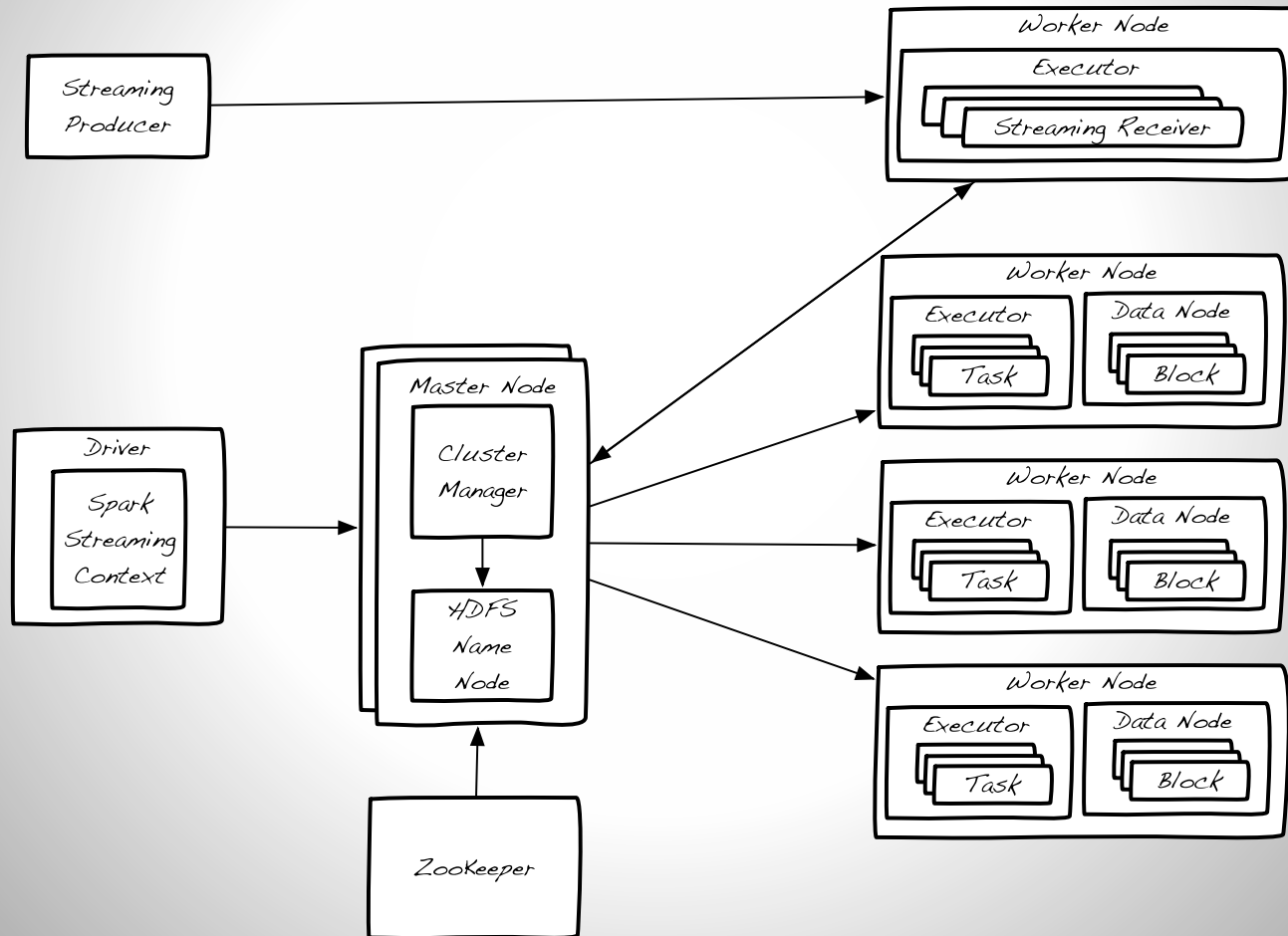
Window Operations

# Similar Projects

- Spark Streaming
  - Twitter Storm
  - Yahoo! S4
  - LinkedIn Samsa
  - Google Millwheel

# Spark Streaming Use Cases

- Operational dashboards
- ETL on streaming data ingestion
- Anomaly, malware, and fraud detection
- Predictive maintenance
  - Sensors
- NLP analysis
  - Twitter Firehose
- Lambda architecture
  - Unified batch and streaming
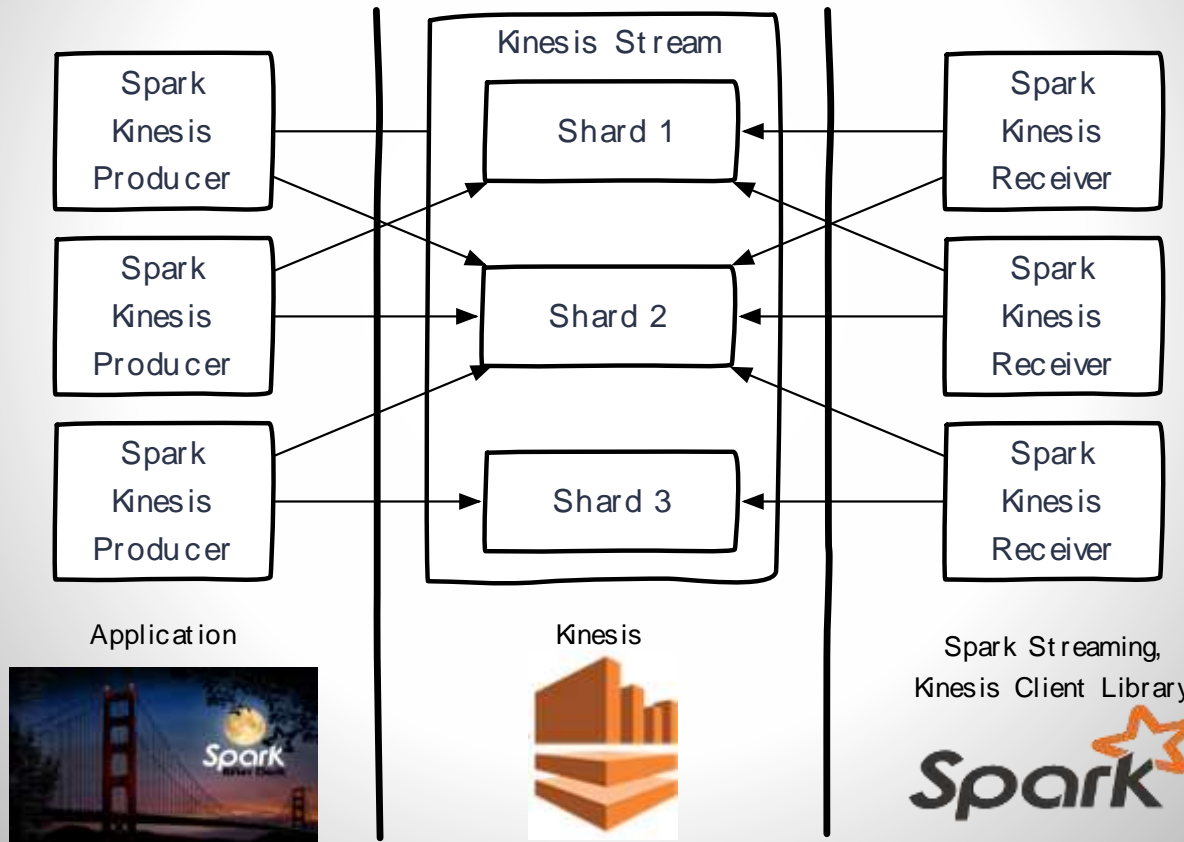  - ie. Different machine learning models for different time frames

# Spark Streaming Cluster View

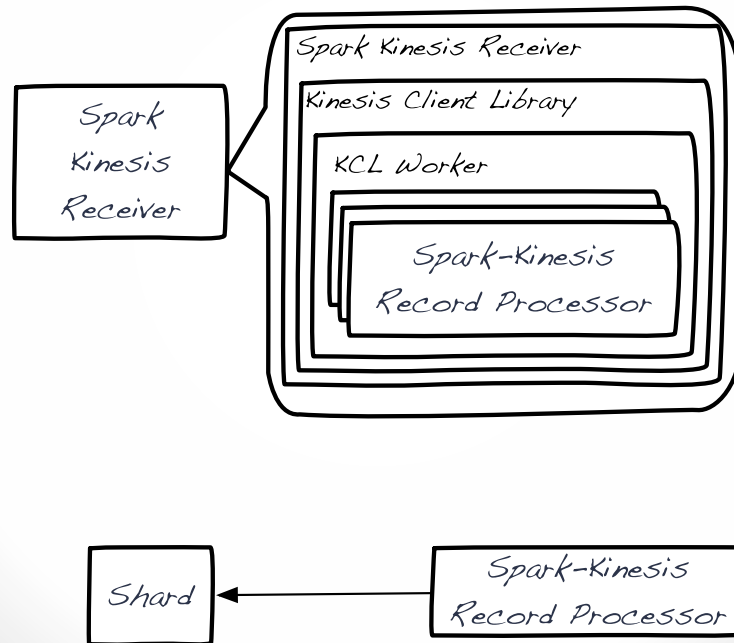Spark Streaming + Hadoop Cluster View

# Spark Streaming + Kinesis



Spark Streaming Kinesis Architecture

# Spark Kinesis Receiver Internals

# Scaling

- Horizontally scale by adding more Kinesis Receivers
- Kinesis Client Library within each Kinesis Receiver will negotiate and rebalance shard processing
- Never need more Kinesis Receivers than the number of stream shards
- Each Kinesis Receiver can support multiple shards
- Supports Kinesis shard splitting/merging
- Recommendation:  over provision shards and avoid splitting/merging

# Demo!



https://github.com/apache/spark/blob/master/extras/kinesis-asl/src/main/...

Scala:  scala/org/apache/spark/examples/streaming/KinesisWordCountASL.scala

Java:  java/org/apache/spark/examples/streaming/JavaKinesisWordCountASL.java

# Streaming Receiver Failure

- Upon failure, backup receiver takes over
- Checkpoint sources like Kafka and kinesis allow multiple receivers to pull from the same stream (ie. during a failover)
  - De-duping is handled by Spark
- Supports graceful shutdown to allow in-flight message draining
- Recommendation:  choose buffered sources like Flume, Kafka and Kinesis

# Streaming Driver Failure

- Streaming Driver app is long-running
  - Monitor driver, receiver worker nodes, and streams
  - Alert upon failure or unusually high latency
- Driver failure interrupts stream processing
- Enable checkpointing for backup Driver to take over
- Use StreamingContext.getOrCreate(…) in Driver app

# Types of Checkpoints

Spark

1. Spark checkpointing of StreamingContext DStreams and metadata

2. Lineage of state and window DStream operations

Kinesis

3. Kinesis Client Library (KCL) checkpoints current position within shard

   – Checkpoint info is stored in DynamoDB per Kinesis Application keyed by shard

# Monitoring

- Streaming tab in Spark Web UI
- CloudWatch
- StreamingListener callback

# Web UI Monitoring

# Tuning

- Batch interval
  - High:  reduce overhead of submitting new tasks for each batch
  - Low:  keeps latencies low
  - Recommendation:  test & find the sweet spot
- Checkpoint interval
  - High:  reduce load on checkpoint overhead
  - Low:  reduce amount of data loss on failure
  - Recommendation:  5-10x sliding window interval
- Explicitly uncache DStreams when no longer needed
- Use CMS GC for consistent processing times
- Use Kryo serialization
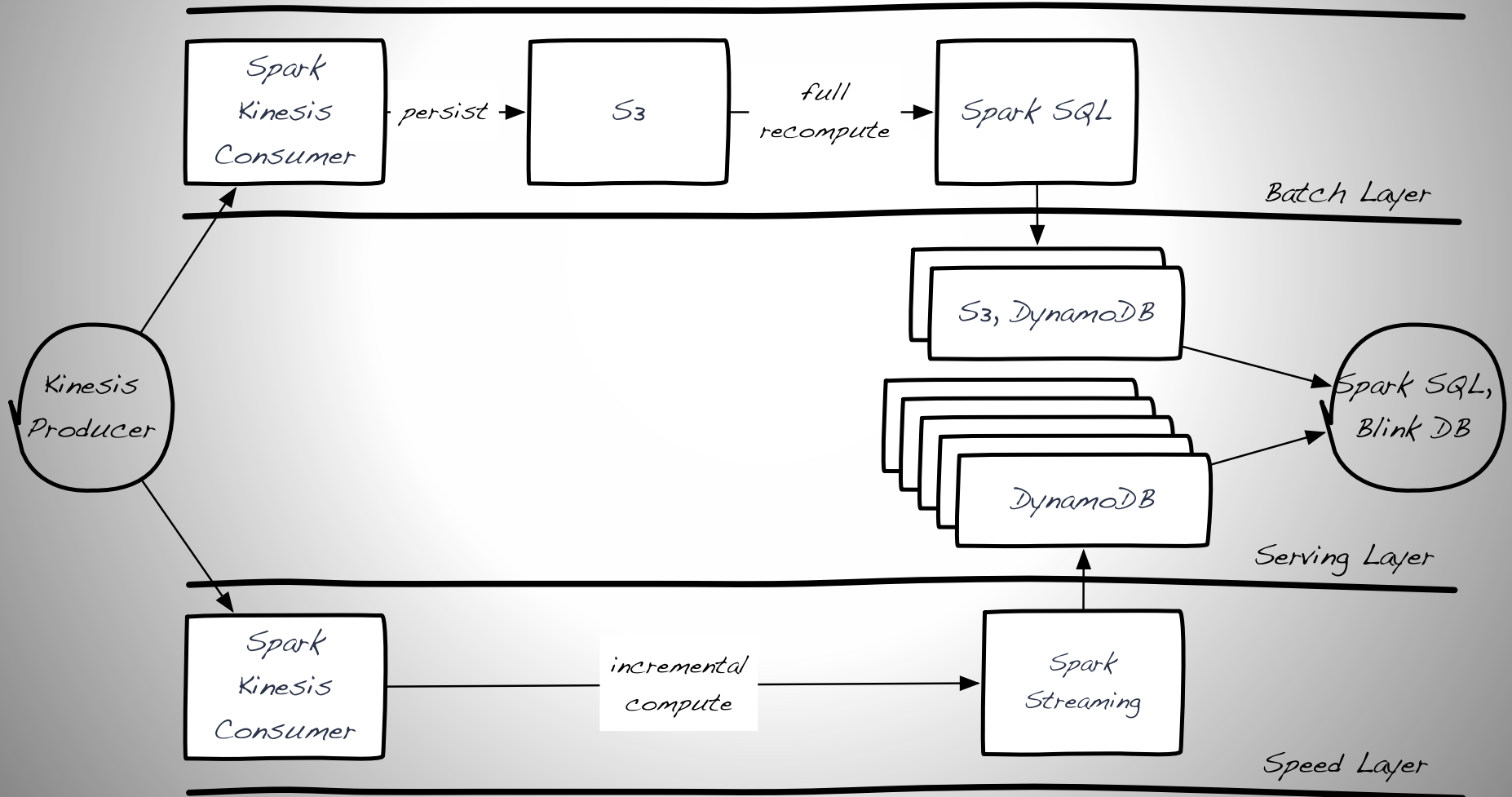- DStreams are already serialized as byte arrays (versus Java objects) to minimize GC

# Lambda Architecture Overview

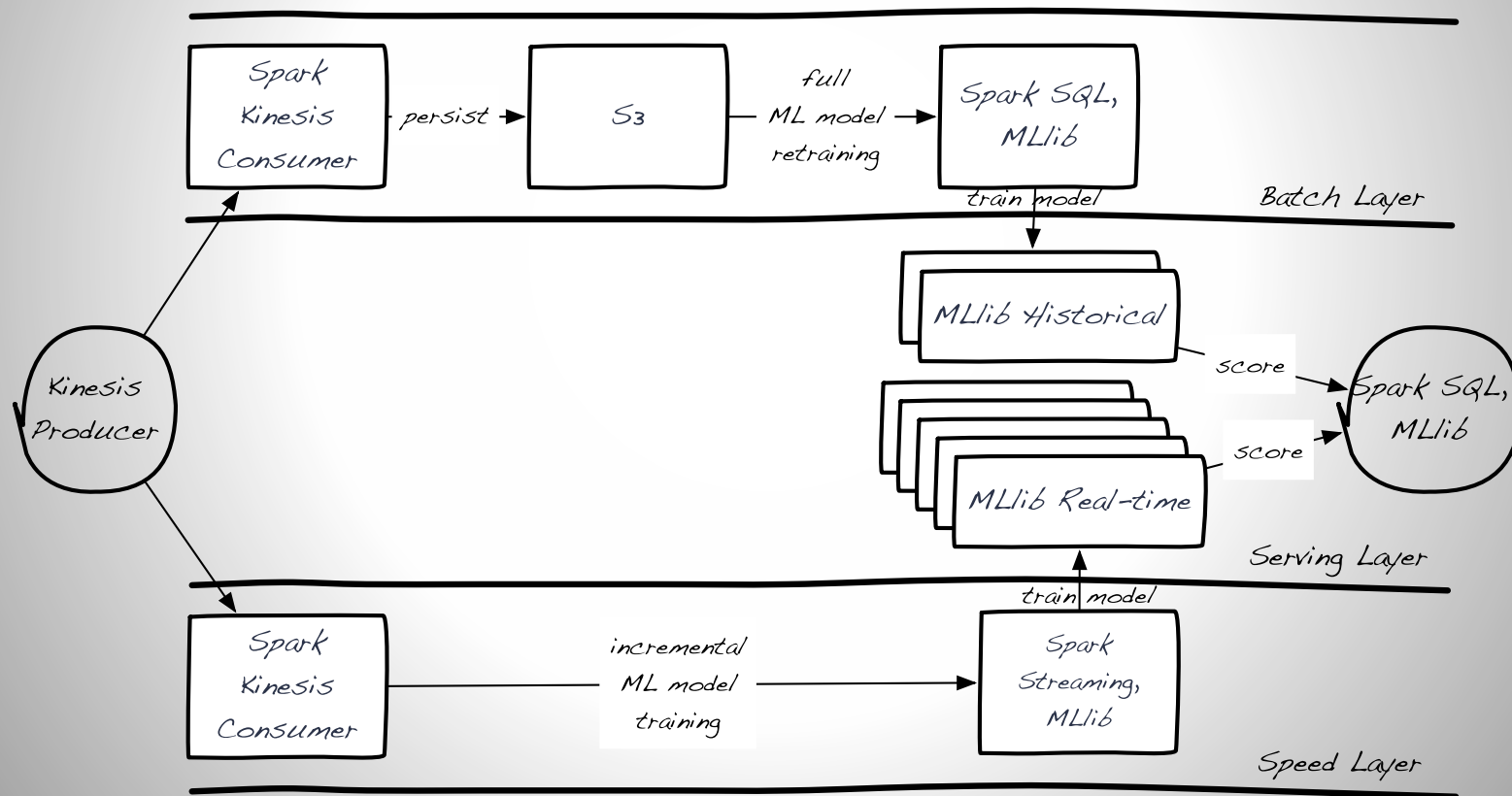Lambda Architecture as Initially Described

# Spark + AWS + Lambda

## Lambda Architecture Spark and AWS

# Spark + AWS + Lambda + ML

Lambda Architecture Spark, AWS, and Machine Learning

# Summary

- Spark
- Spark Streaming + Kinesis
- Scaling
- Fault Tolerance
- Monitoring
- Tuning
- Lambda Architecture
- Spark in Action
  - Oct 2014 MEAP
  - Early access:  http://sparkinaction.com