# Data Pipelines

cs5356
Daniel Doubrovkine
@dblockdotorg

# *Moving Data Around*



Compute

Stateless
Ephémère
Non résilient

Storage

Statefull
Persistant
Résilient

# *Typical Warehousing: Batch*

CORNELL NYC TECH

**Event Data**
billions of data points
terabytes of data

→

Extract
Transform
Load

→

**Insight**

# *Too slow … continuous processing*

**Event Stream**
billions of events
terabytes of data

Join
Organize
Aggregate

**Insight**

# *Clickstream*

- What users do, usually in log files.
  date, time, IP, URL

- Page views per URL over time?
- Top N page views of all time?
- What products to visitors buy together?

# *Aggregate*

- Distributed Messaging System
  - **In:** Logs
  - **Out:** Domain Data

kafka

# *Hadoop*

- ## Storage: HDFS

  Linux ext3

  Replicated blocks

  Write Once / Read Many

  Retry transfers

- ## Execution: Map/Reduce

  In parallel on many servers

  Retry on failure

  Map + Reduce

  Servers can join or leave

# *Map/Reduce Canonical Example*

- Word Count
  - **Map**: transform text to { word: 1 }
  - **Reduce**: sum by key

# *Replacing Parts of Hadoop*

- Spark
- Microsoft Dryad
- Apache Tez
- Impala
- Google Big Query
- Google Cloud Dataflow

- ## Richer API
  filter(), join(), distinct(), groupByKey()

- ## Maintain State
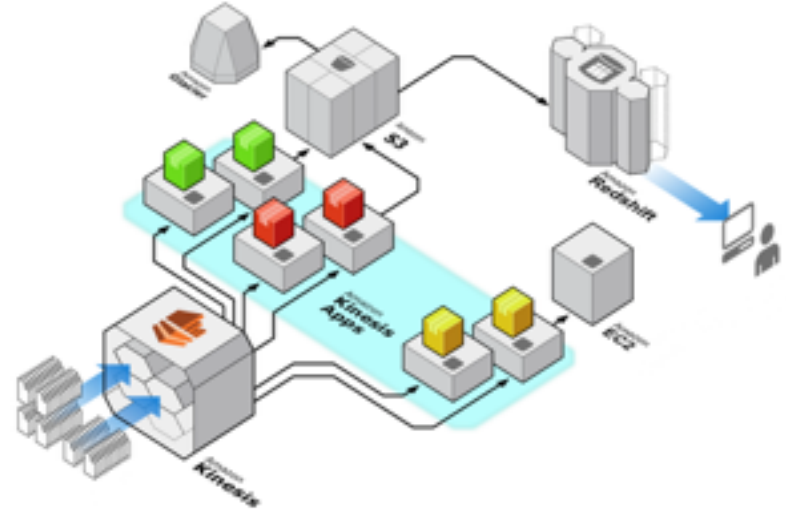  updateStateByKey()
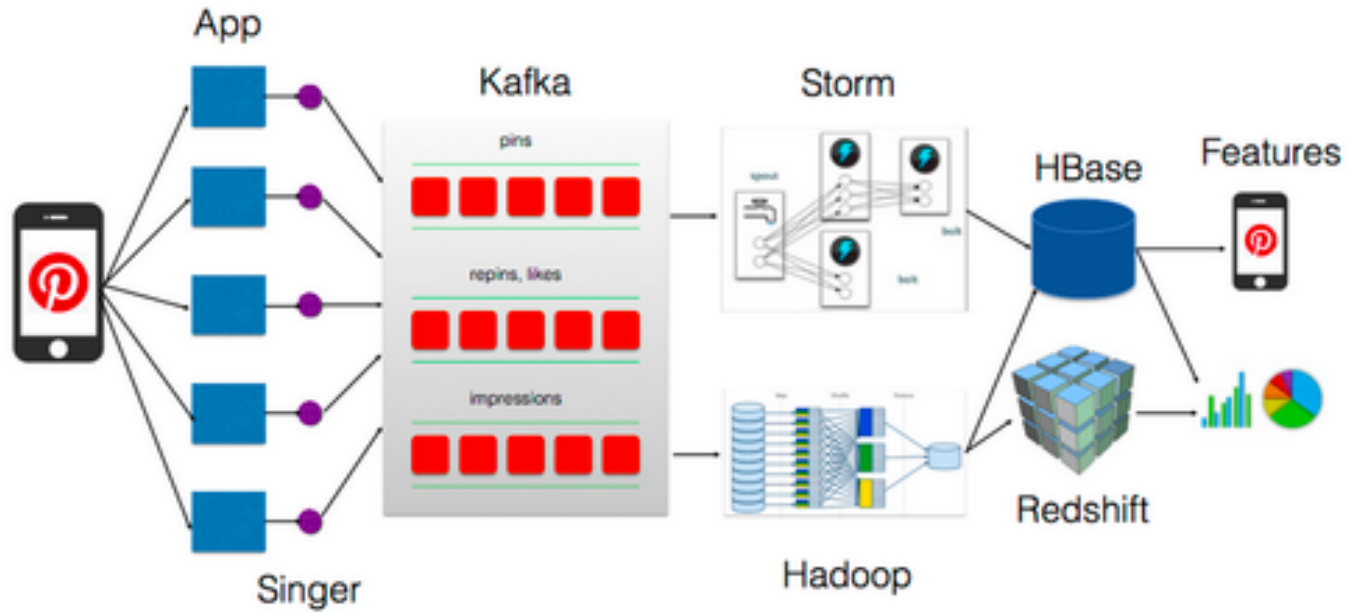
- ## Window Operations
  window size and slide interval

# *Amazon Kinesis*



- Easy
- Real-Time
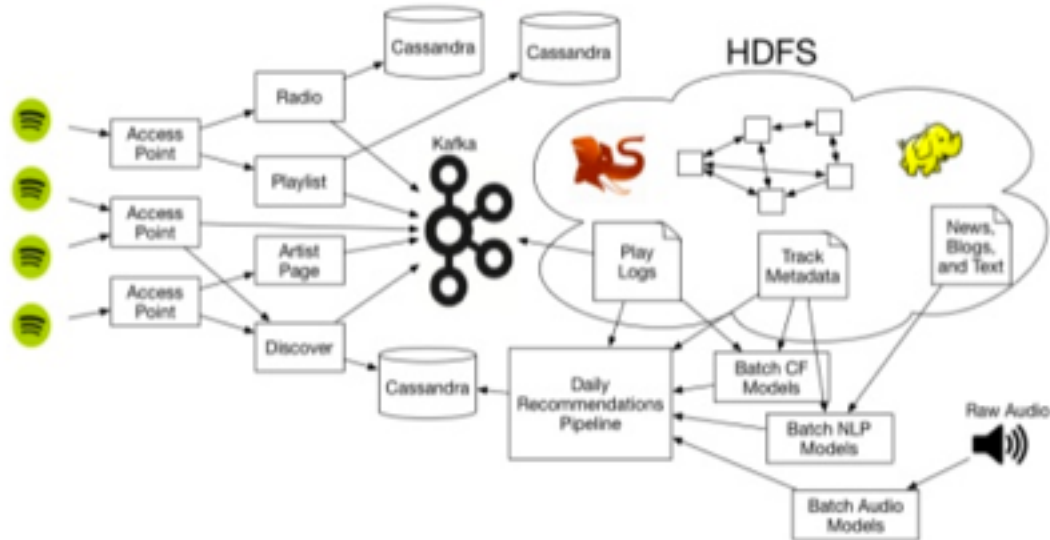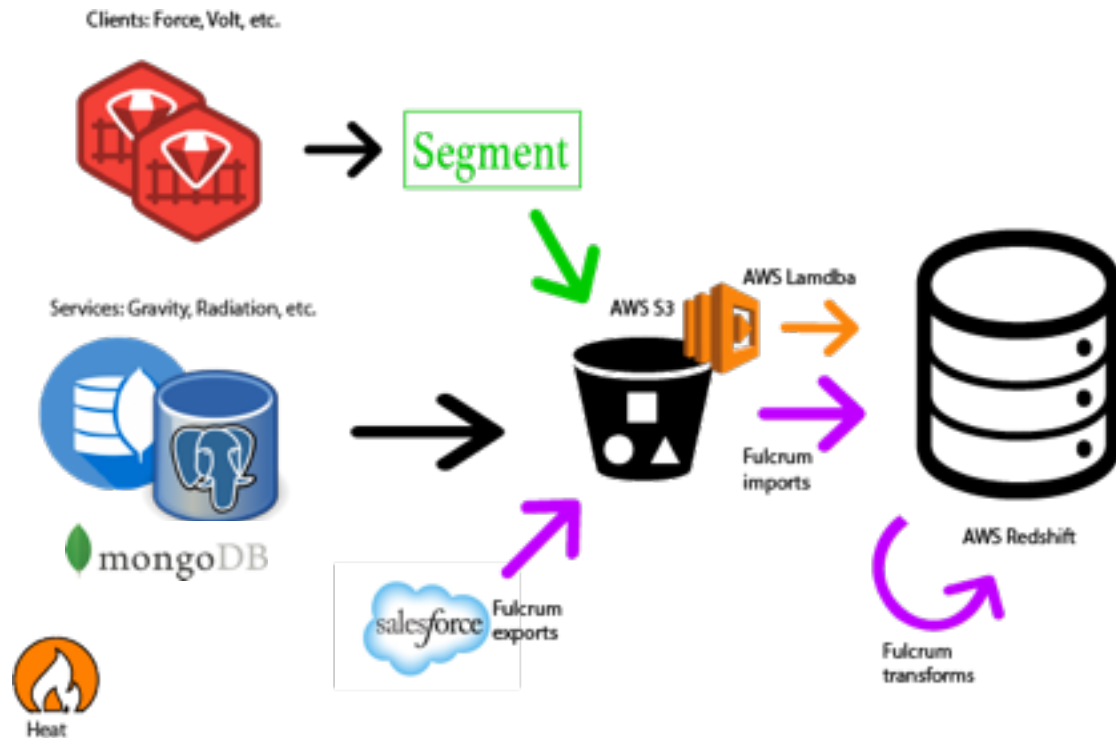- High Throughput
- S3/DynamoDB/Redshift

# *Pinterest*



Data Architecture overview

# *Spotify*

# *Artsy*

# MongoDB (after the break)

Shannon Bradshaw