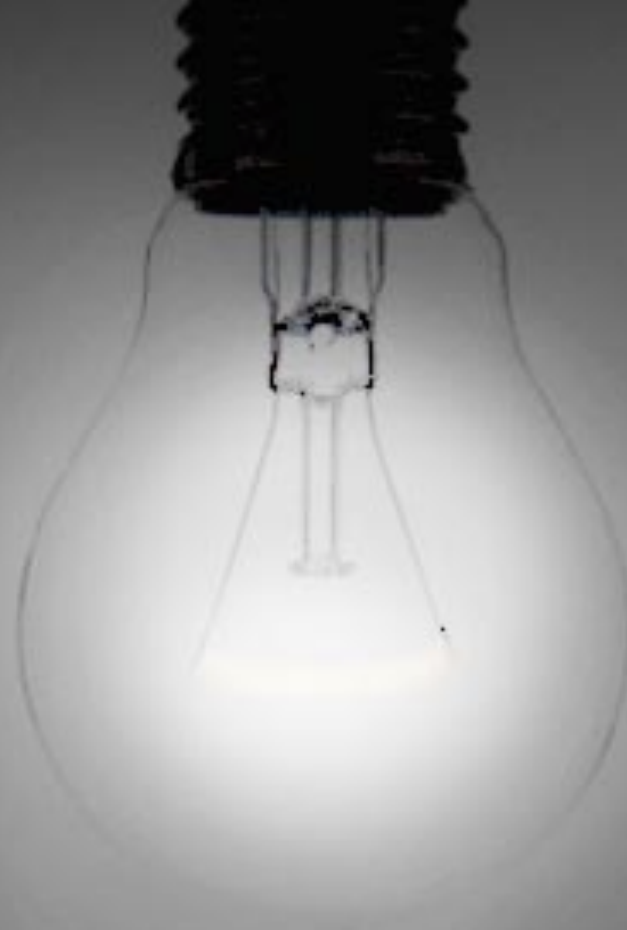Typesafe

# You say big data
# I say fast data

Nilanjan Raychaudhuri
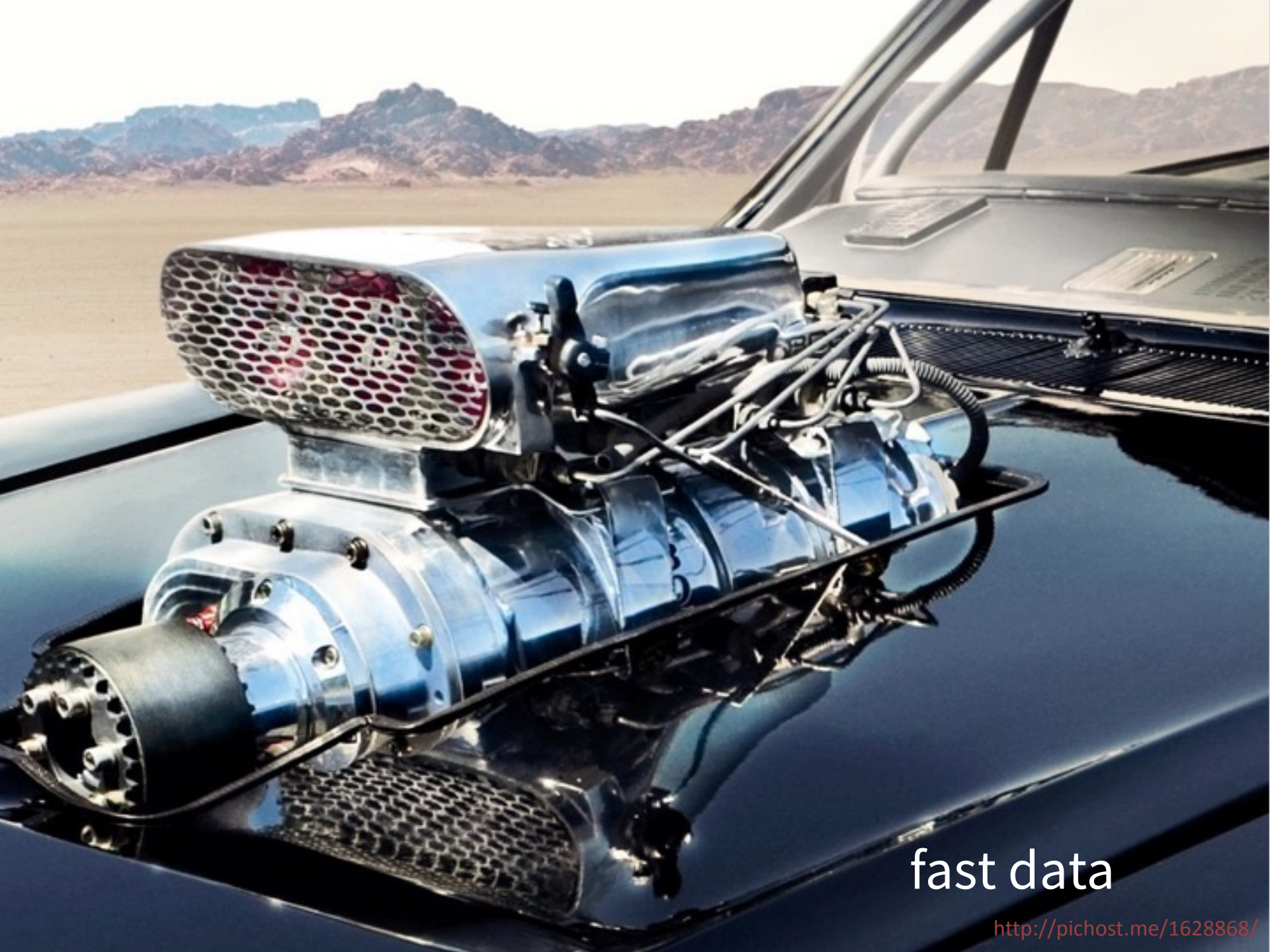
@nraychaudhuri

Big data

# Big Data

Word for storing, managing and making money of very large dataset
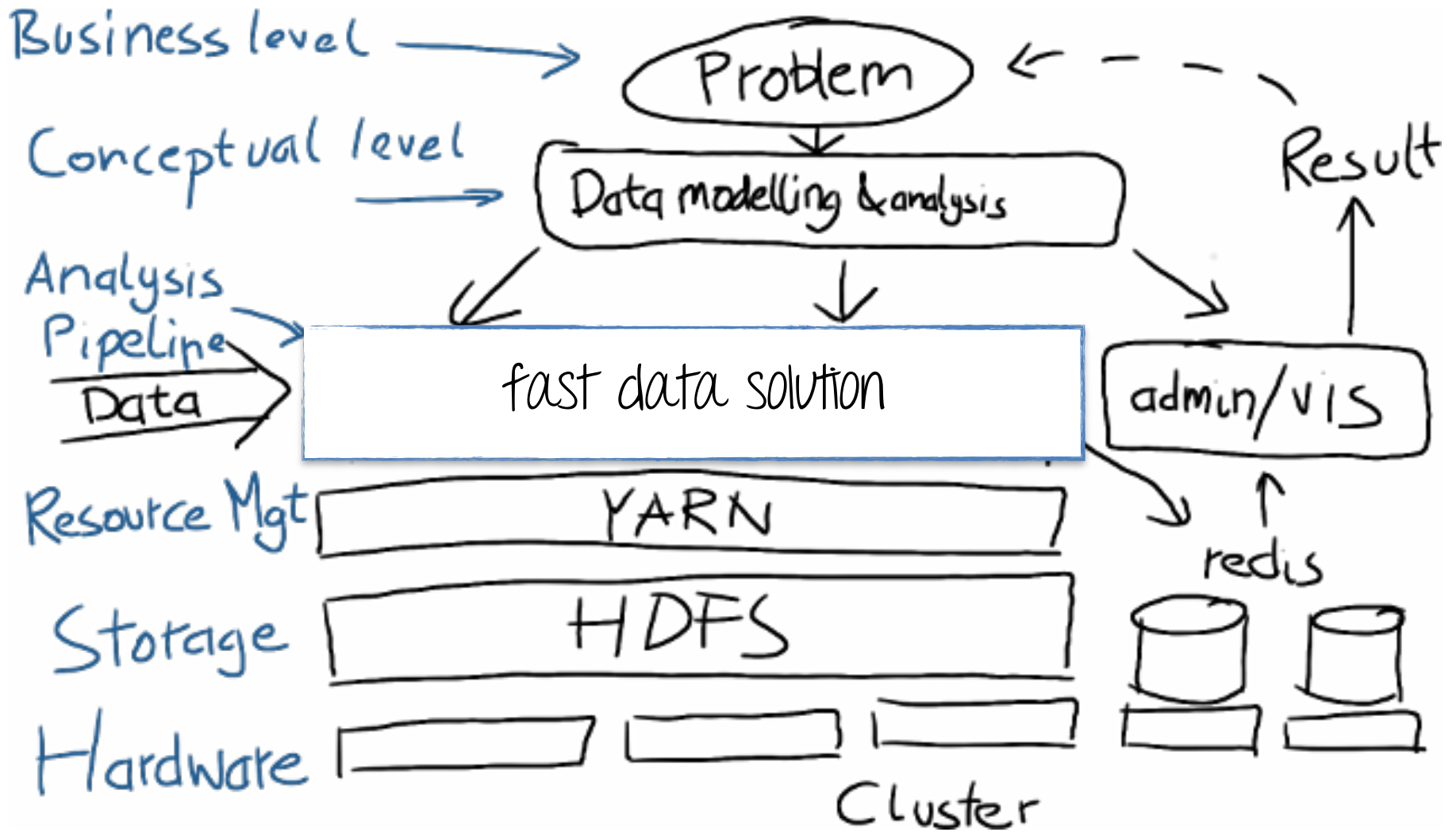
fast data

# Fast Data

Word for using fast, real time analytics, online machine learning for profit
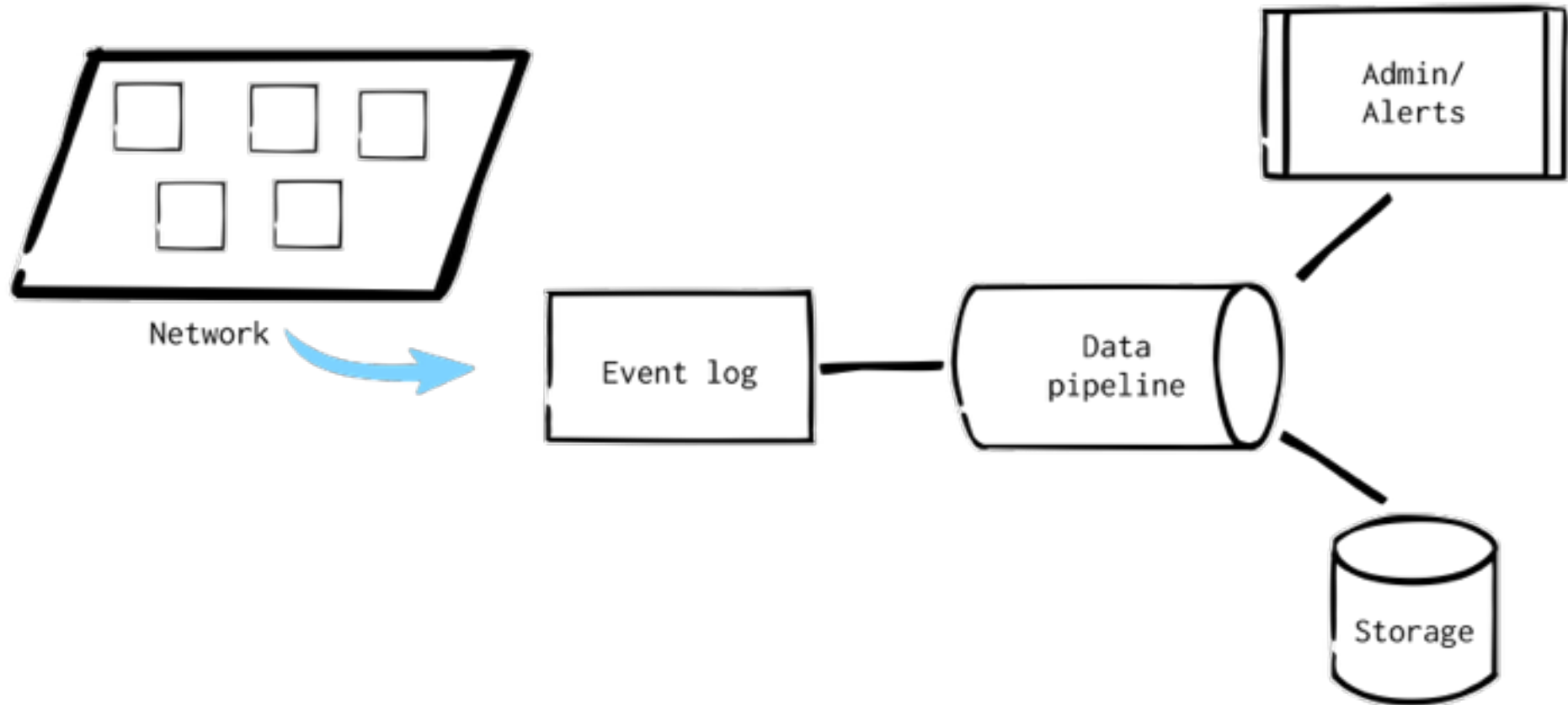
# Definition

The phrase ***Fast Data*** captures this range of new systems and approaches, which balance various tradeoffs to deliver timely, cost-efficient streaming data processing.

"You can't expect the value of data to just appear out of thin air. Data isn't fissile material. It doesn't spontaneously reach critical mass and start producing insights."

- Marko Karppinen

# Detecting network intrusion

Network

Event log

Data pipeline

Admin/ Alerts

Storage

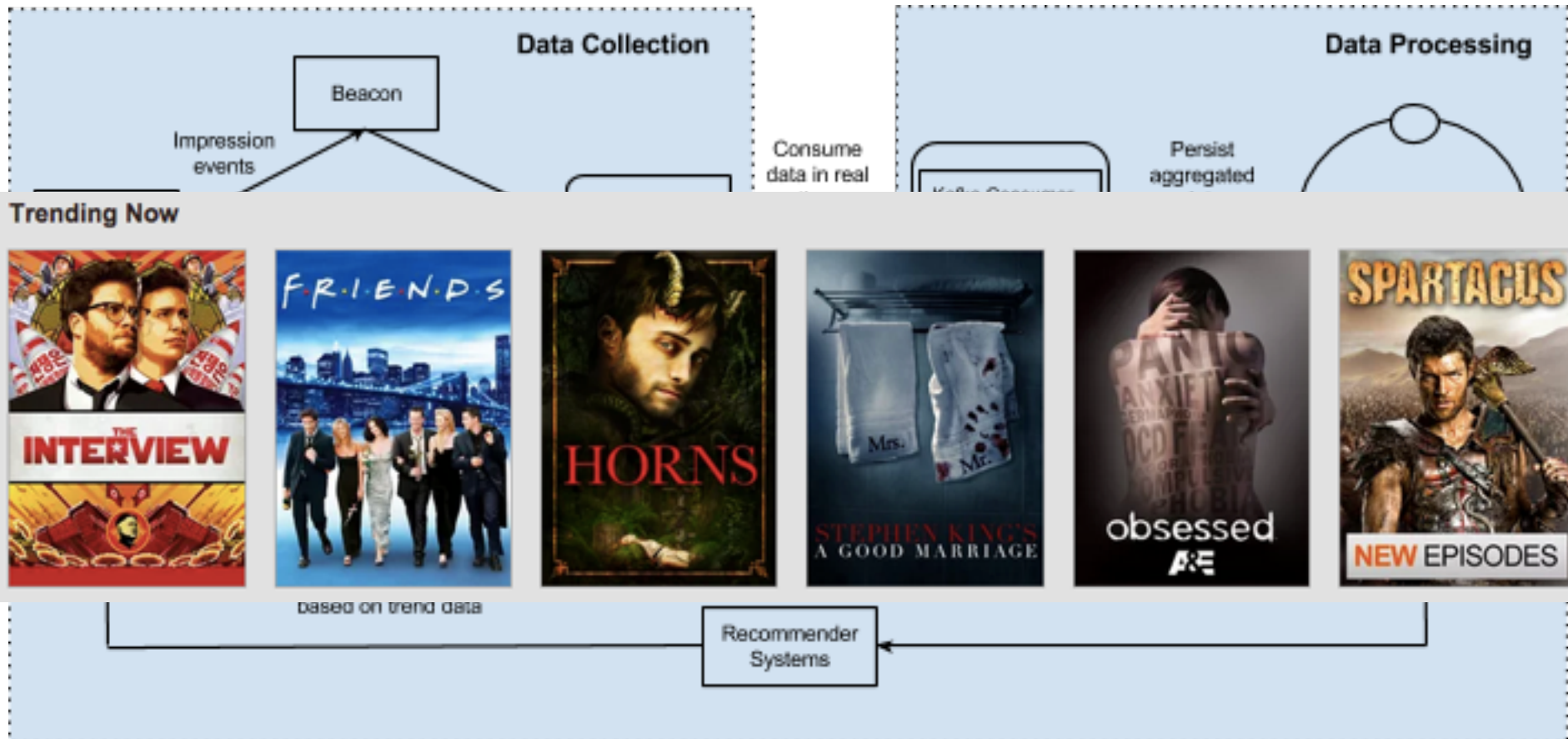# Faster data analysis @ Pinterest
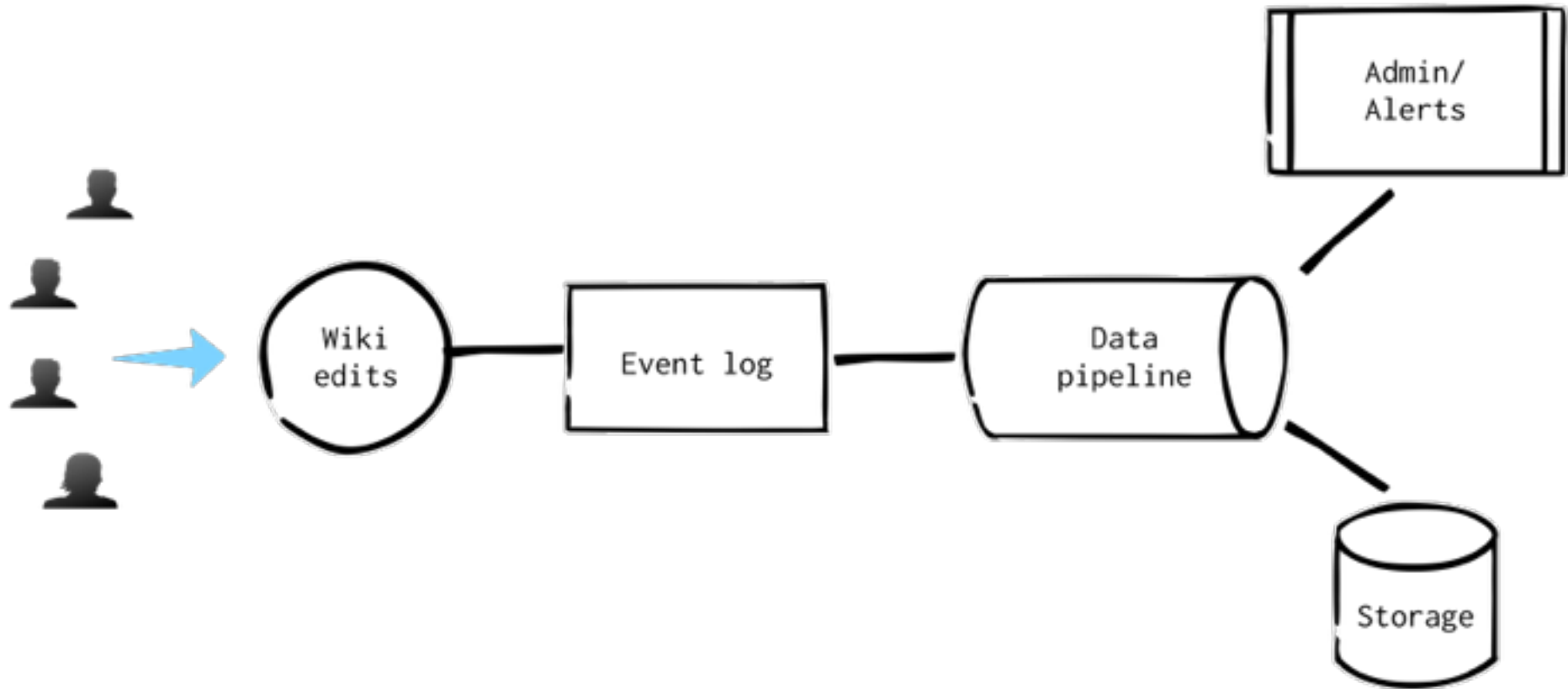


http://engineering.pinterest.com/post/111380432054/real-time-analytics-at-pinterest

# Recommendation engine @ Netflix



http://techblog.netflix.com/2015/02/whats-trending-on-netflix.html
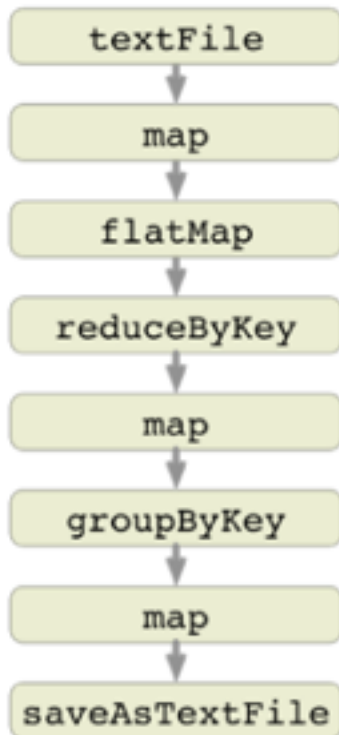
# Predict breaking news

What big data/fast data has to do with
<u>functional programming</u>?

# Functional programming is the killer Paradigm for Data Apps

- Dean Wampler, Ph.D

# Dataflow programming

```
         textFile

            map

          flatMap

        reduceByKey

            map

         groupByKey

            map

       saveAsTextFile
```

```scala
sparkContext.textFile("/path/to/input")
.map { line =>
  val array = line.split(",", 2)
  (array(0), array(1))
}.flatMap {
  case (id, contents) => toWords(contents).map(w => ((w,id),1))
}.reduceByKey {
  (count1, count2) => count1 + count2
}.map {
  case ((word, path), n) => (word, (path, n))}
.groupByKey
.map {
  case (word, list) => (word, sortByCount(list))
}.saveAsTextFile("/path/to/output")
```
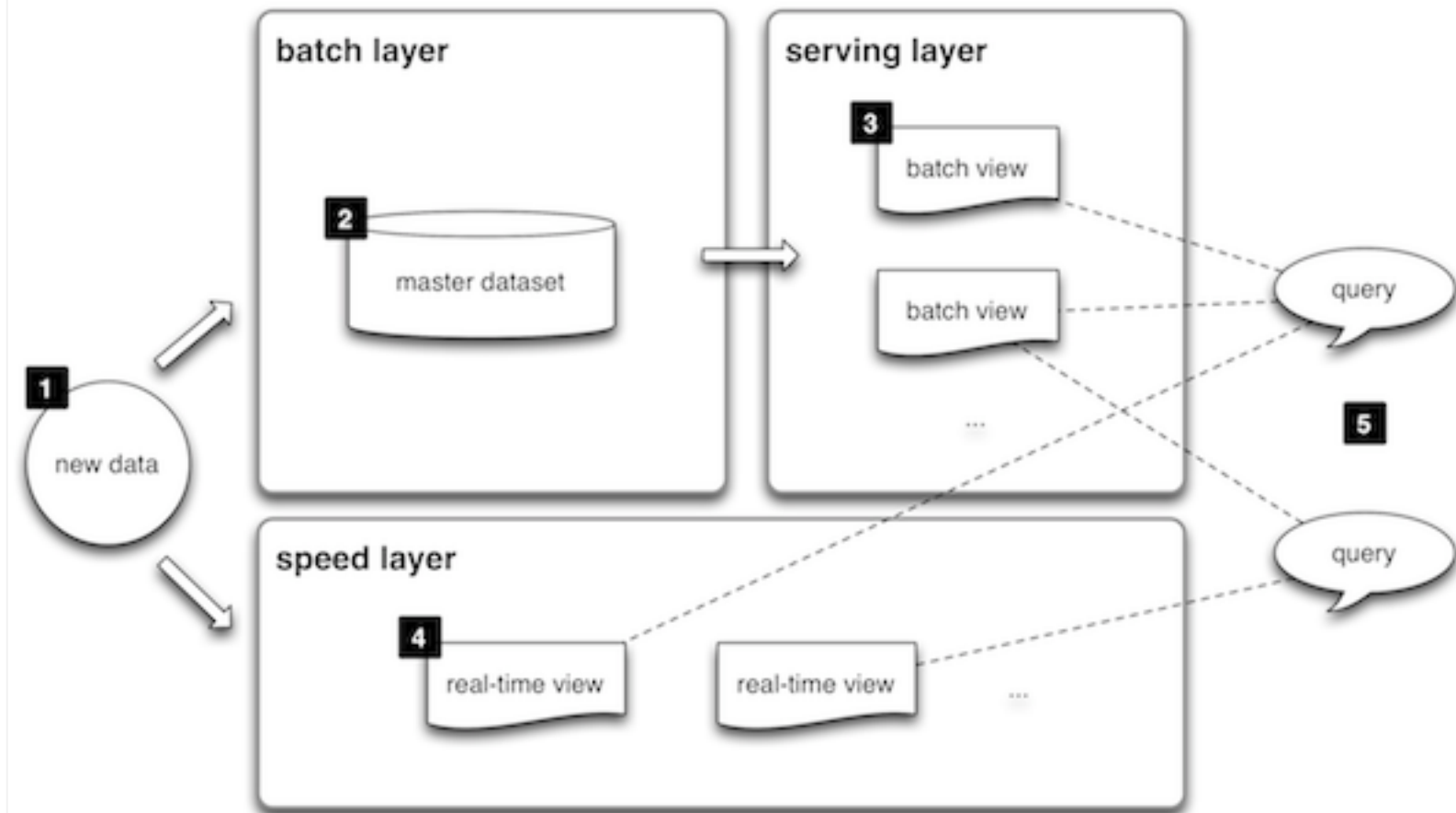
# What is streaming?

# Define: Streaming

- A type of data processing engine that is designed with infinite data sets in mind

- Other common terms for streaming: unbounded data, low latency, approximate, and/or speculative results

# Hurdles

# Lambda architecture

# Typical Lambda architecture

# Latency

# new tools…

## Highly Scalable Blog

Articles on Big Data, NoSQL, and Highly Scalable Software Engineering

## Probabilistic Data Structures for Web Analytics and Data Mining

Posted on May 1, 2012

https://highlyscalable.wordpress.com/2012/05/01/probabilistic-structures-web-analytics-data-mining/

# We need new tools

SELECT avg(sessionTime)
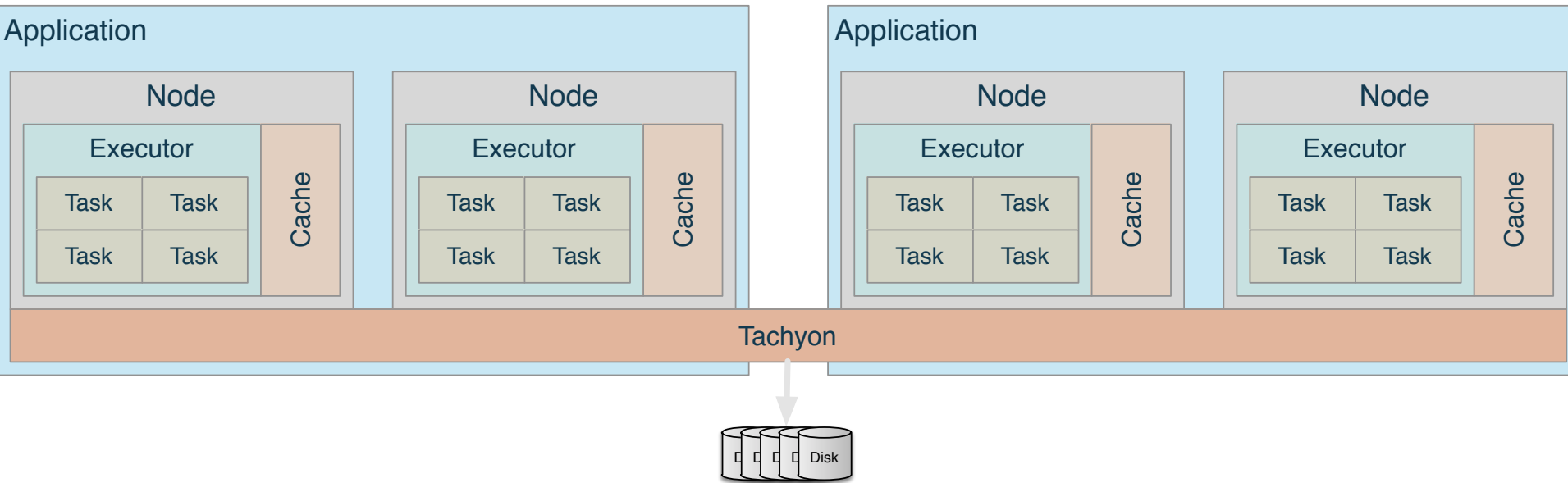FROM Table
WHERE city='San Francisco'
WITHIN 2 SECONDS

Queries with Time Bounds

SELECT avg(sessionTime)
FROM Table
WHERE city='San Francisco'
ERROR 0.1 CONFIDENCE 95.0%

Queries with Error Bounds

BlinkDB

# new tools…

| Application | | | | Application | | | |
|---|---|---|---|---|---|---|---|
| **Node** | | | | **Node** | | | |
| Executor | | Cache | | Executor | | Cache | |
| Task | Task | | | Task | Task | | |
| Task | Task | | | Task | Task | | |

**Tachyon**

Disk

Typesafe

# Online ML
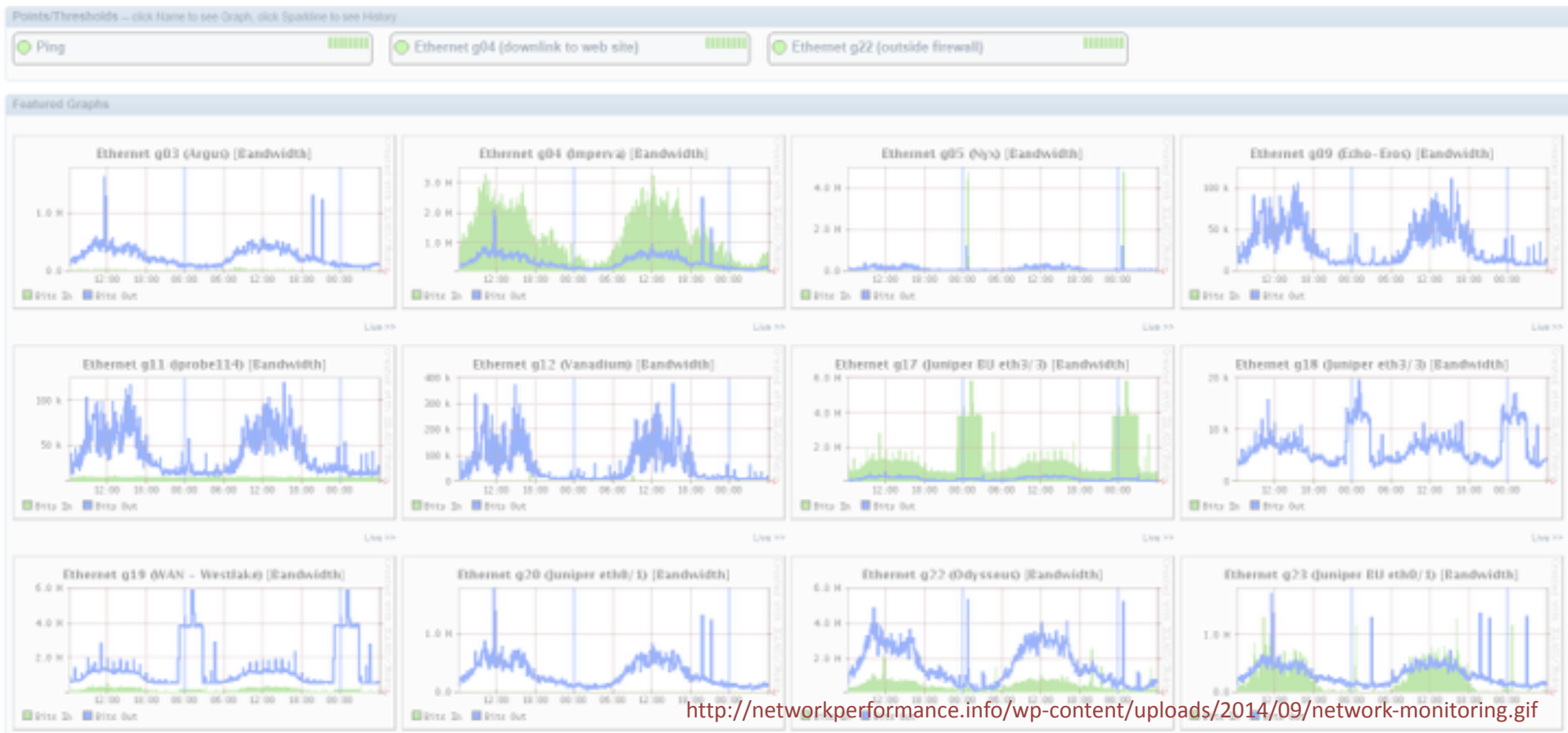
# Computation model

```java
public static class LineIndexMapper
 extends MapReduceBase
 implements Mapper<LongWritable, Text,
                   Text, Text> {
 private final static Text word =
  new Text();
 private final static Text location =
  new Text();

 public void map(
  LongWritable key, Text val,
  OutputCollector<Text, Text> output,
  Reporter reporter) throws IOException {

  FileSplit fileSplit =
   (FileSplit)reporter.getInputSplit();
```

# Monitoring



http://networkperformance.info/wp-content/uploads/2014/09/network-monitoring.gif

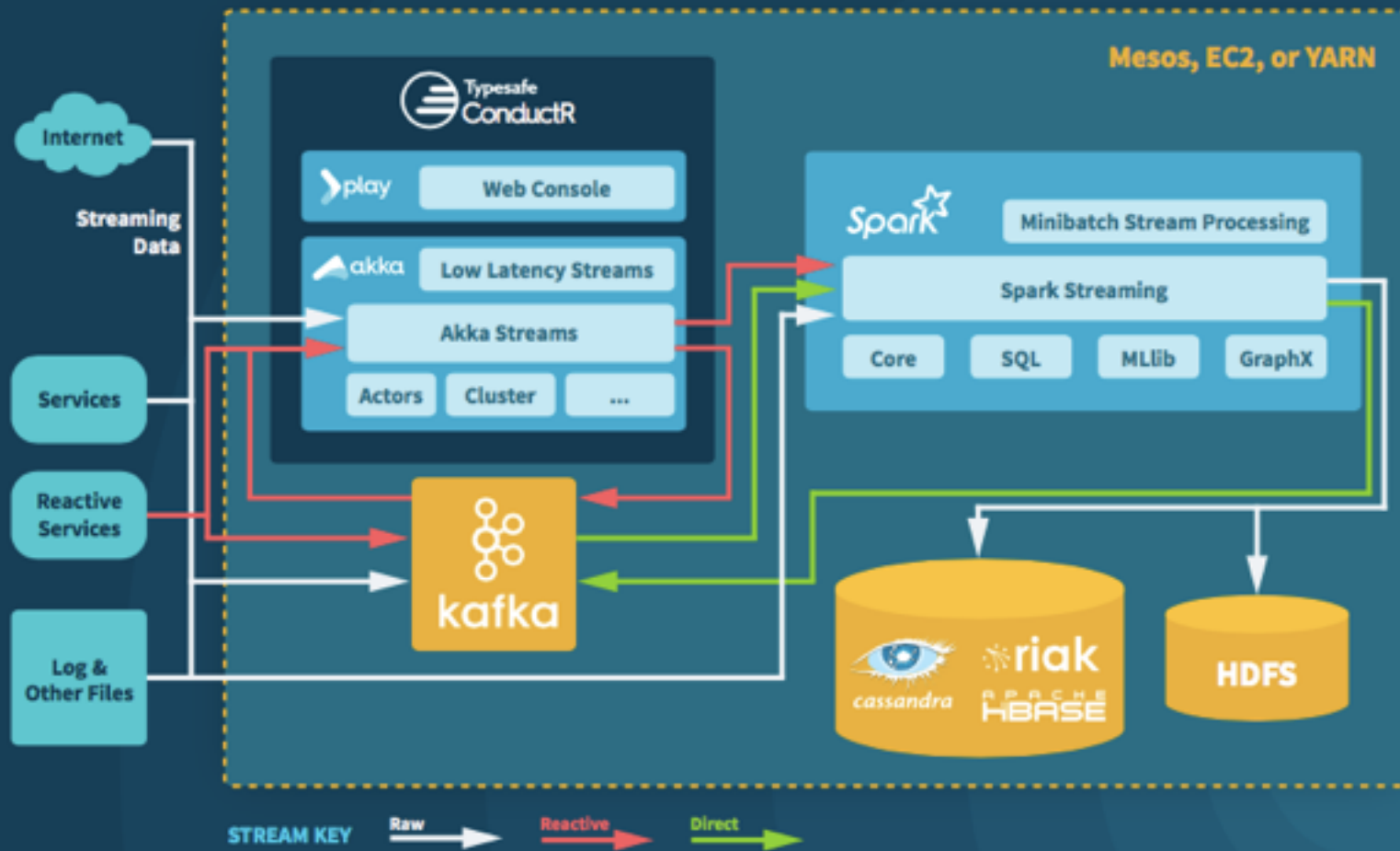# Wish list

# Streaming platform feature list

- Should be programmer friendly (Scala and FP)

- Should be fault tolerant

- Should support high throughput/low latency

- Should integrate with batch systems (Hadoop)

- Should run machine learning algorithm

-  Should not overwhelm the consumer

- Should provide integration points for other streaming systems

Typesafe

And that streaming platform is…
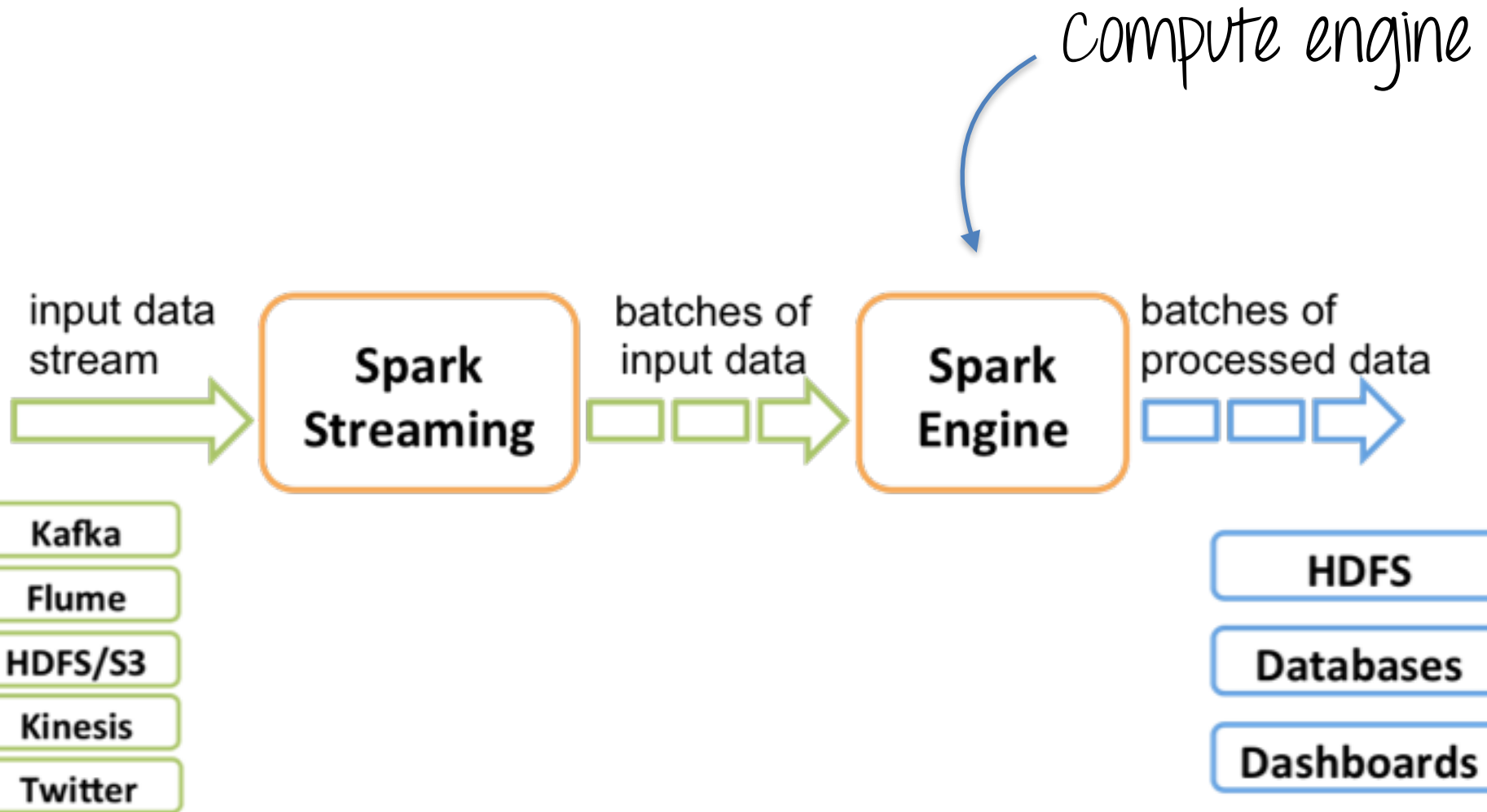
not there yet

but...something is emerging

# Fast data platform

# Spark Streaming

*Captures time slices of events*

# Data pipeline

*Compute engine*

input data stream

**Spark Streaming**

batches of input data

**Spark Engine**

batches of processed data

Kafka
Flume
HDFS/S3
Kinesis
Twitter

HDFS
Databases
Dashboards

Typesafe

Show me the code

# Predict breaking news

# Under the hood

Spark Driver ("main")

```
sc = SparkContext
ssc = StreamingContext(
  sc, Seconds(5))
```

Cluster Manager

**Node**

Executor

| Task | Task |
|------|------|
| Task | Task |

Cache

**Node**

Executor

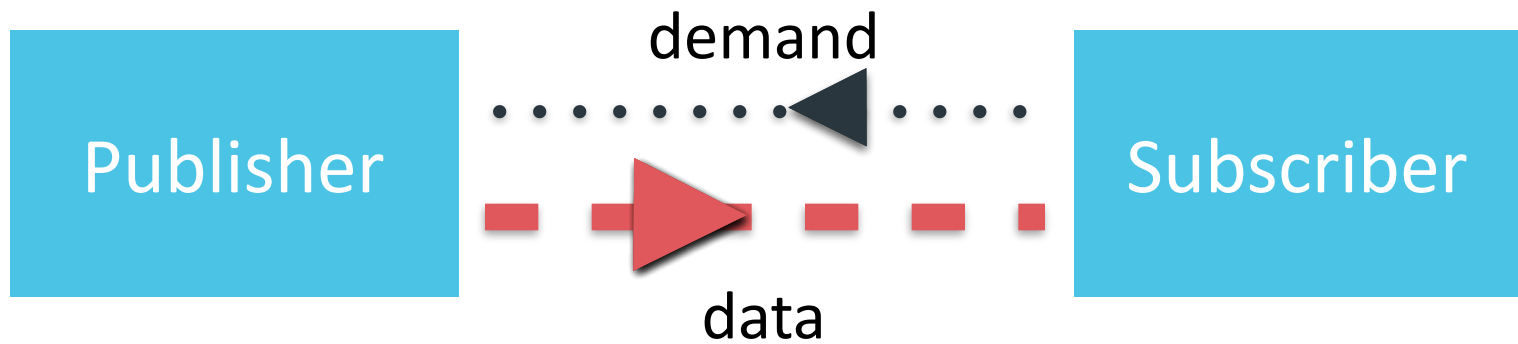| Task | Task |
|------|------|
| Task | Task |

Cache

...

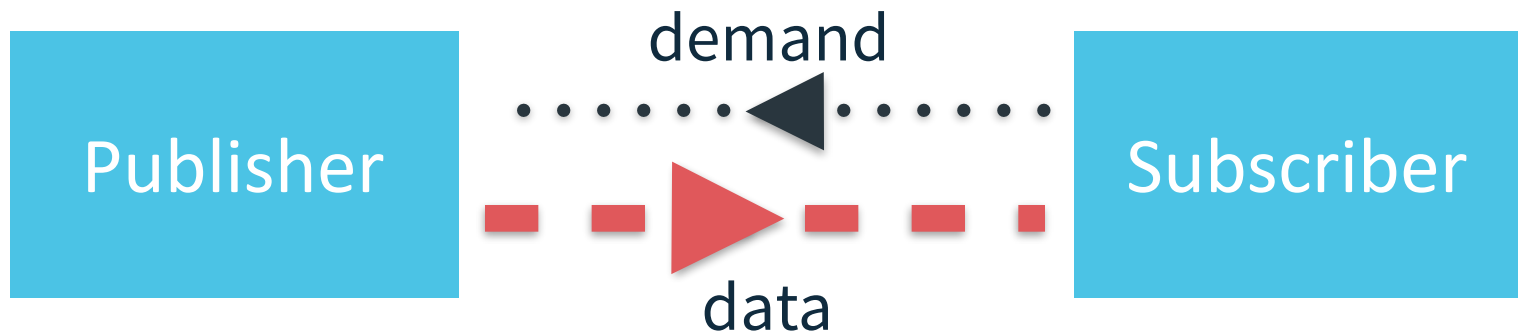*Reactive Streams*

# Definition

A standard for asynchronous stream processing with non-blocking back pressure

# Supply and Demand

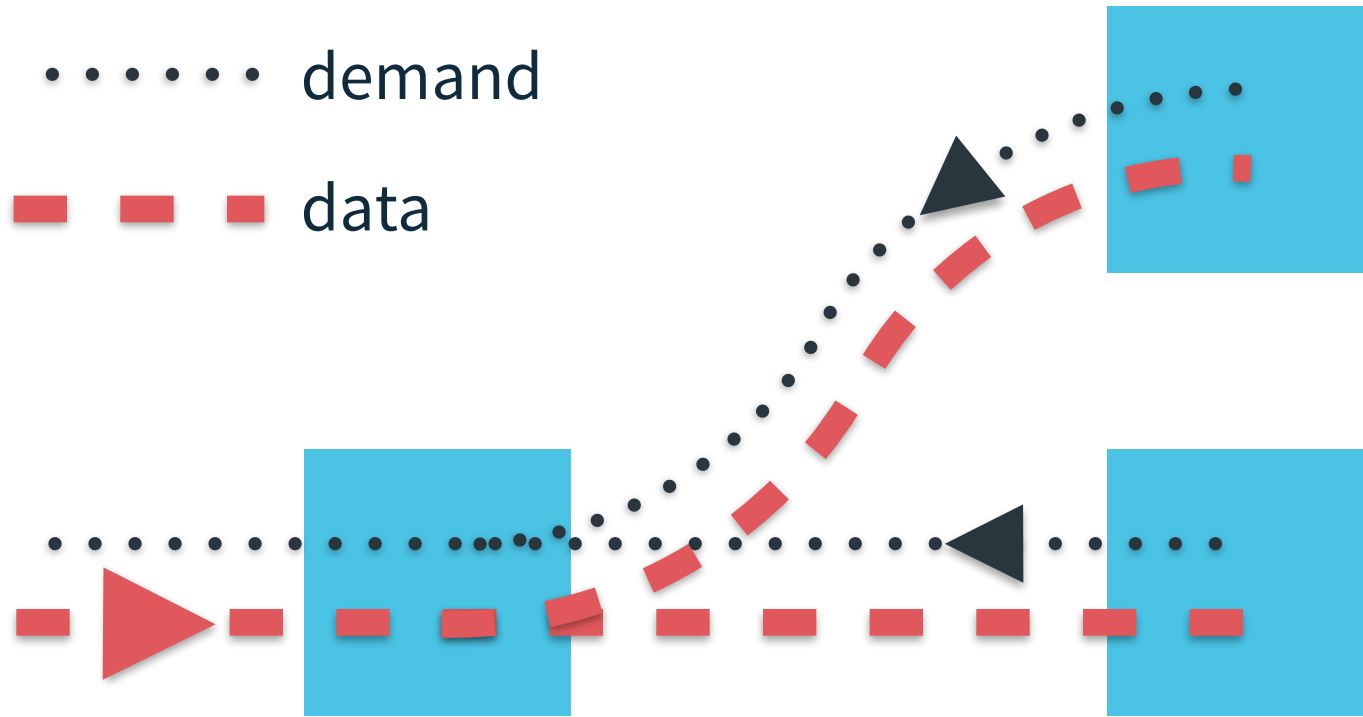# Dynamic Push–Pull

- "push" behavior when consumer is faster
- "pull" behavior when producer is faster

Publisher
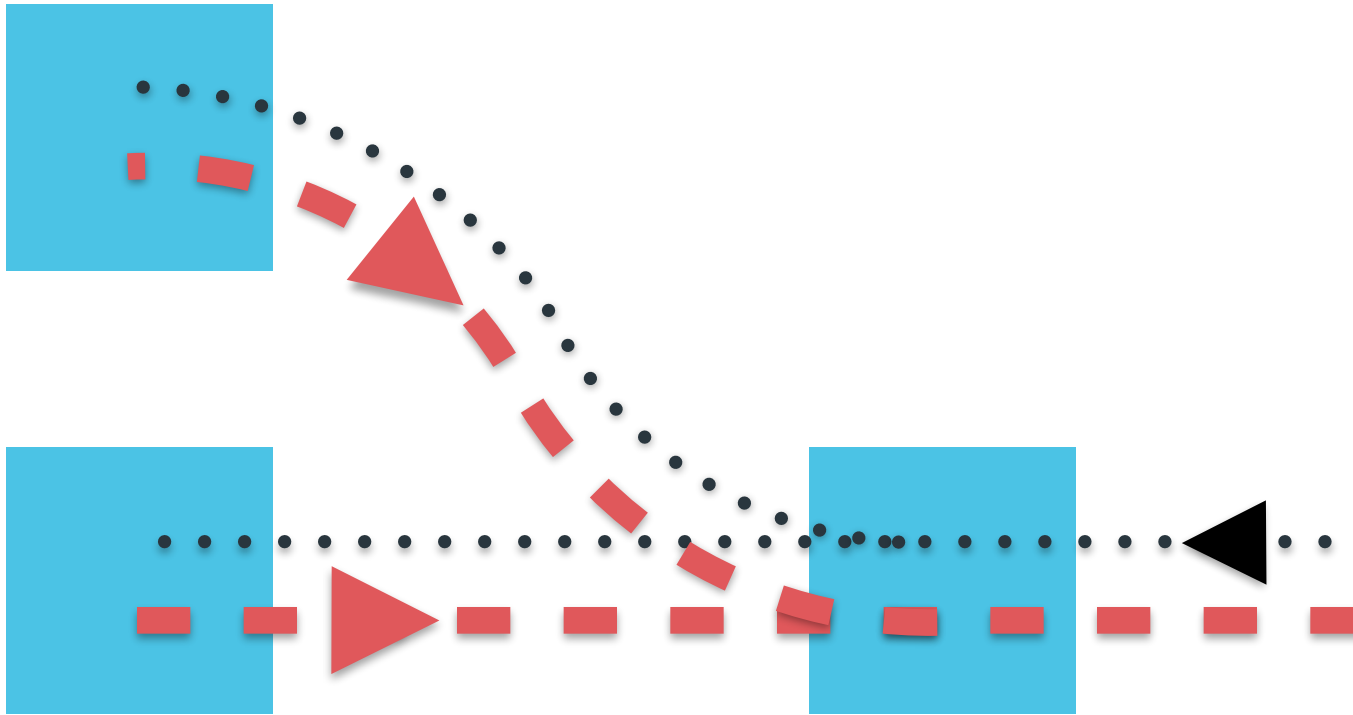
demand

data

Subscriber

**Typesafe**

# Explicit Demand: Tailored Flow Control

······ demand

▬ ▬ ▬ ▬ data

splitting the data means merging the demand

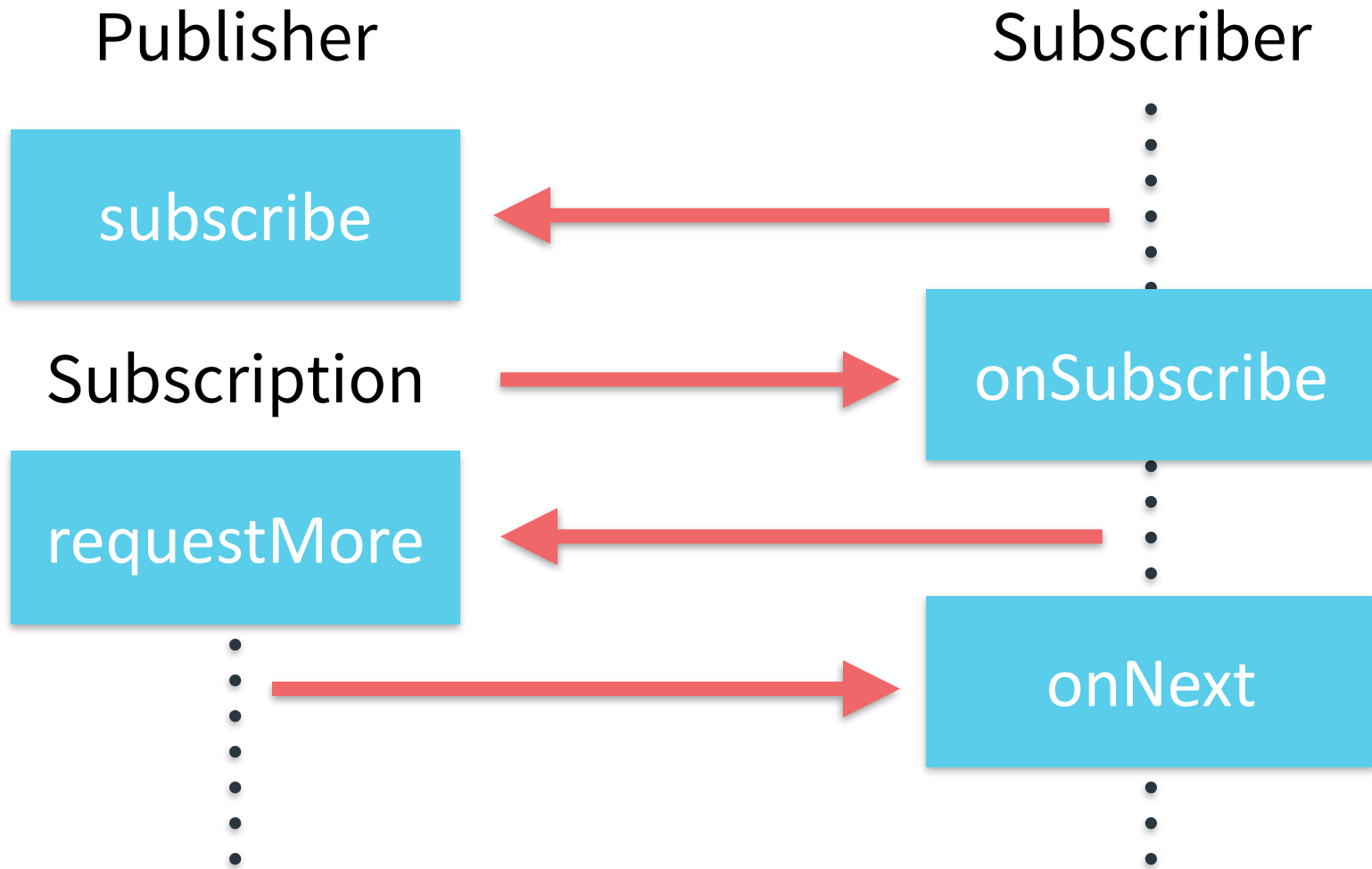# Explicit Demand: Tailored Flow Control



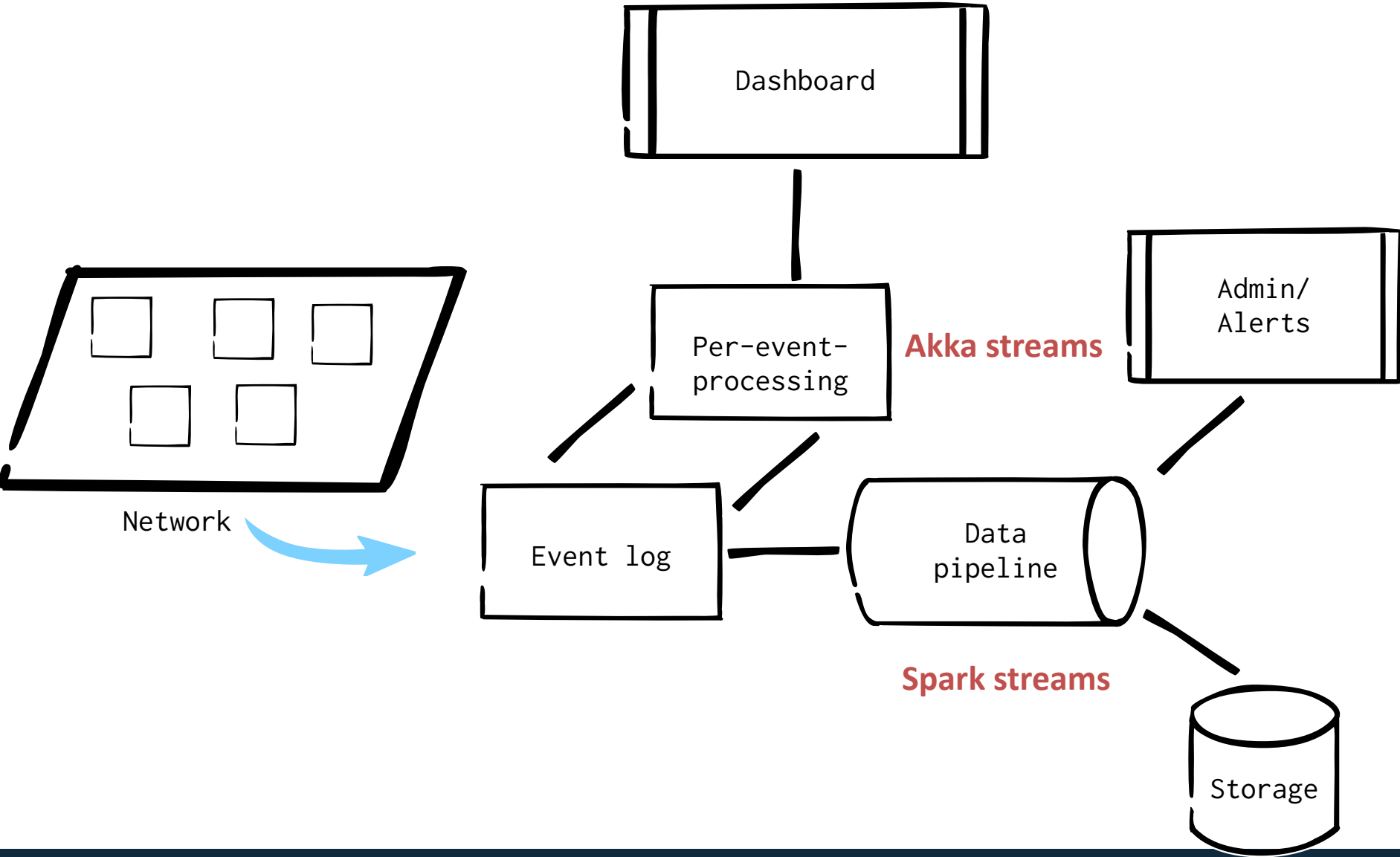merging the data means splitting the demand

# The Meat

```scala
trait Publisher[T] {
  def subscribe(sub: Subscriber[T]): Unit
}
trait Subscription {
  def requestMore(n: Int): Unit
  def cancel(): Unit
}
trait Subscriber[T] {
  def onSubscribe(s: Subscription): Unit
  def onNext(elem: T): Unit
  def onError(thr: Throwable): Unit
  def onComplete(): Unit
}
```

# How does it Connect?

Publisher                                    Subscriber

subscribe ←————————————————————

Subscription ——————————→ onSubscribe

requestMore ←————————————————————

——————————→ onNext

Typesafe

# Network intrusion

# Detecting network intrusion

Dashboard

Per-event-processing

**Akka streams**

Admin/
Alerts

Network

Event log

Data
pipeline

**Spark streams**

Storage

Typesafe

```scala
val sc = new SparkContext("local[3]", "Intro")
val ssc = new StreamingContext(sc, Seconds(1))
ssc.checkpoint("data/checkpoint")

val kmeans = new StreamingKMeans()
  .setK(10)
  .setRandomCenters(3, 100.0)
  .setHalfLife(5, "batches")
```

```scala
val receiver = SparkEnv.get.actorSystem.actorOf(trainingDataReceiver,
  "training-data-source")
val trainingStream =
  ssc.actorStream[Vector](TrainingDataStream.props(receiver, kmeans),
  "training-stream")

val rawData: DStream[(String, Vector)] = KafkaUtils
  .createStream(ssc = ssc, kafkaParams, topics,
    StorageLevel.MEMORY_ONLY)
  .map(_._2).map(line => (line, labelAndVector(line)._2))
```

```scala
kmeans.trainOn(trainingStream.map(_._2))
kmeans.predictOnValues(networkStream).print()

ssc.start()
ssc.awaitTermination()
```

# Q & Option[A]

@nraychaudhuri