# Streaming analytics better than batch - when and why ?

**Adam Kawa  -  Dawid Wysakowicz**

getindata

# About Us

- **At GetInData, we build custom Big Data solutions**
  - Hadoop, **Flink**, Spark, Kafka and more
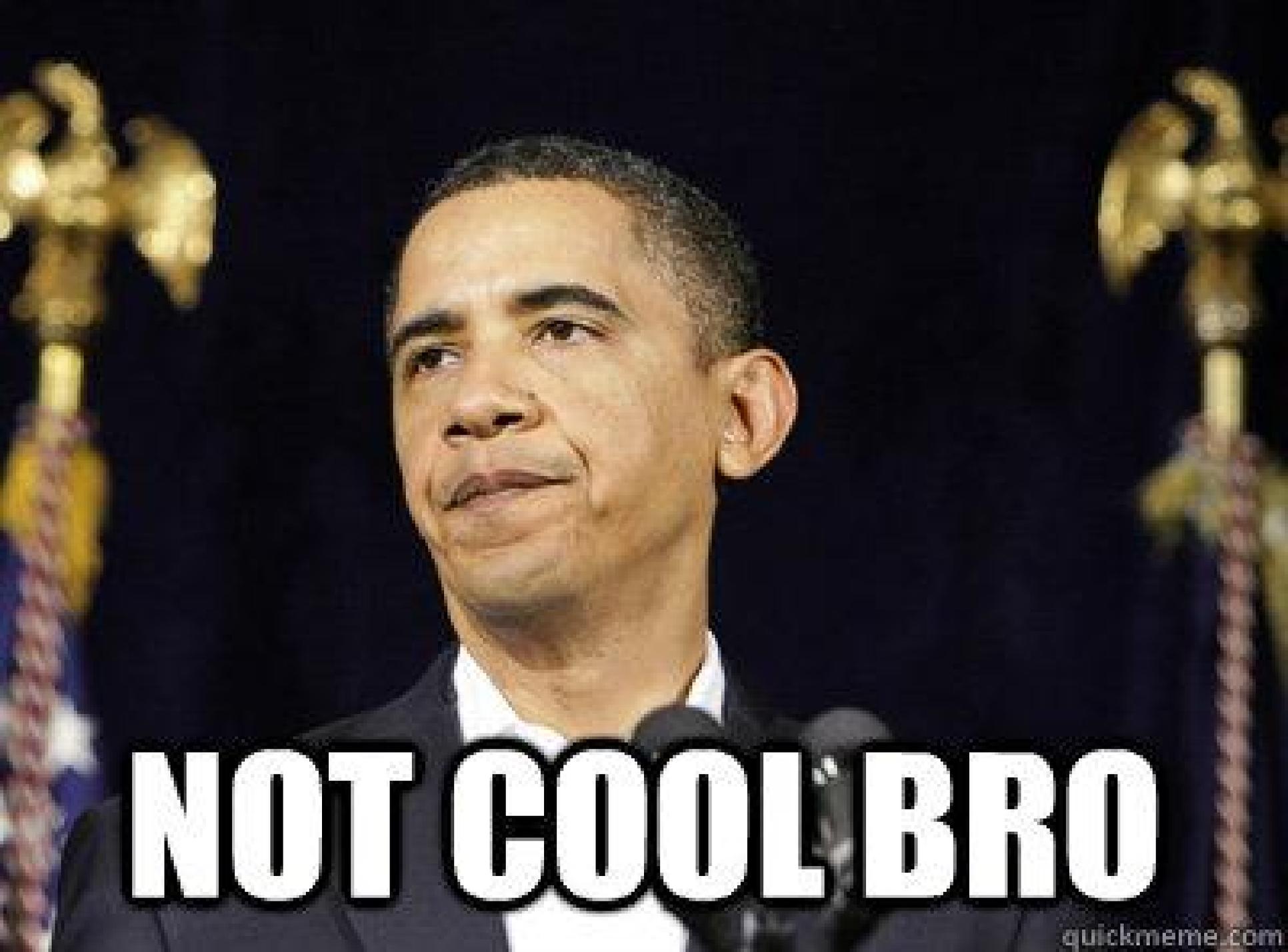- **Our team is today represented by**



**Adam Kawa**



**Dawid Wysakowicz**

Have you ever built cool Big Data pipelines?

NOT COOL BRO

quickmeme.com

# Example Use-Case

- **Can be done in <span style="color:red">batch</span> and <span style="color:red">real-time</span>**
- **User session analytics at Spotify**
  - Simple stats
    - Duration, number of songs, skips, searches etc.
  - Advanced analytics
    - Mood, physical activity, real-time content, ads

Spotify®

# Example Output

**1. Dashboards**



**How long** do users listen to a **new** edition of Discover Weekly?

# Example Output

## 1. Dashboards

## 2. Alerts

**How long** do users listen to a **new** edition of Discover Weekly?

Australian users are listening to Discover Weekly **too short** !!!

# Example Output

## 1. Dashboards



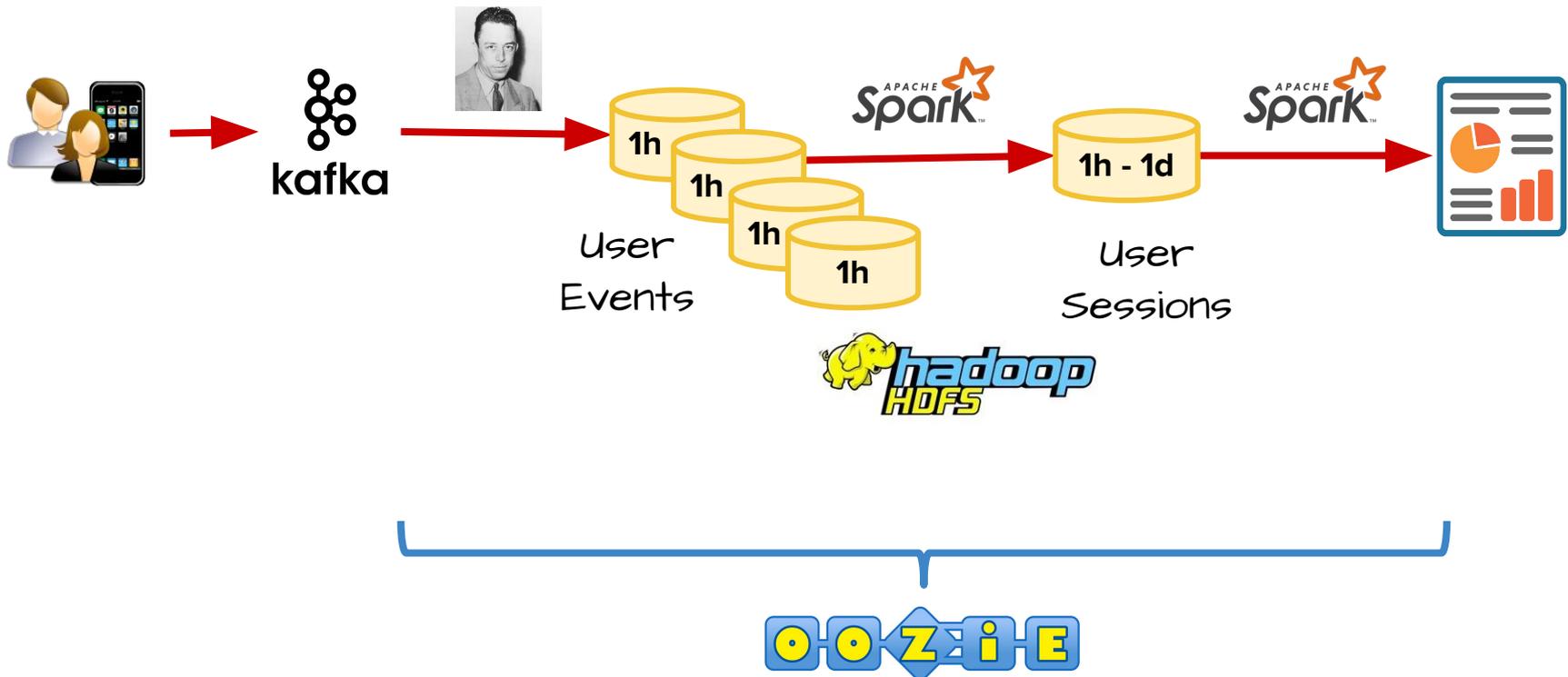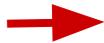**How long** do users listen to a **new** edition of Discover Weekly?

## 2. Alerts



Australian users are listening to Discover Weekly **too short** !!!

## 3. Content



Your Daily Mix 1

**Recommend** songs and ads based on **current** activity.

# 1ˢᵗ - Batch Architecture

# 1st - Batch



FIND OUT MORE!

# Many Moving Parts

⬇ **The higher learning curve**

⬇ **The more gluing code**

⬇ **The larger administrative effort**
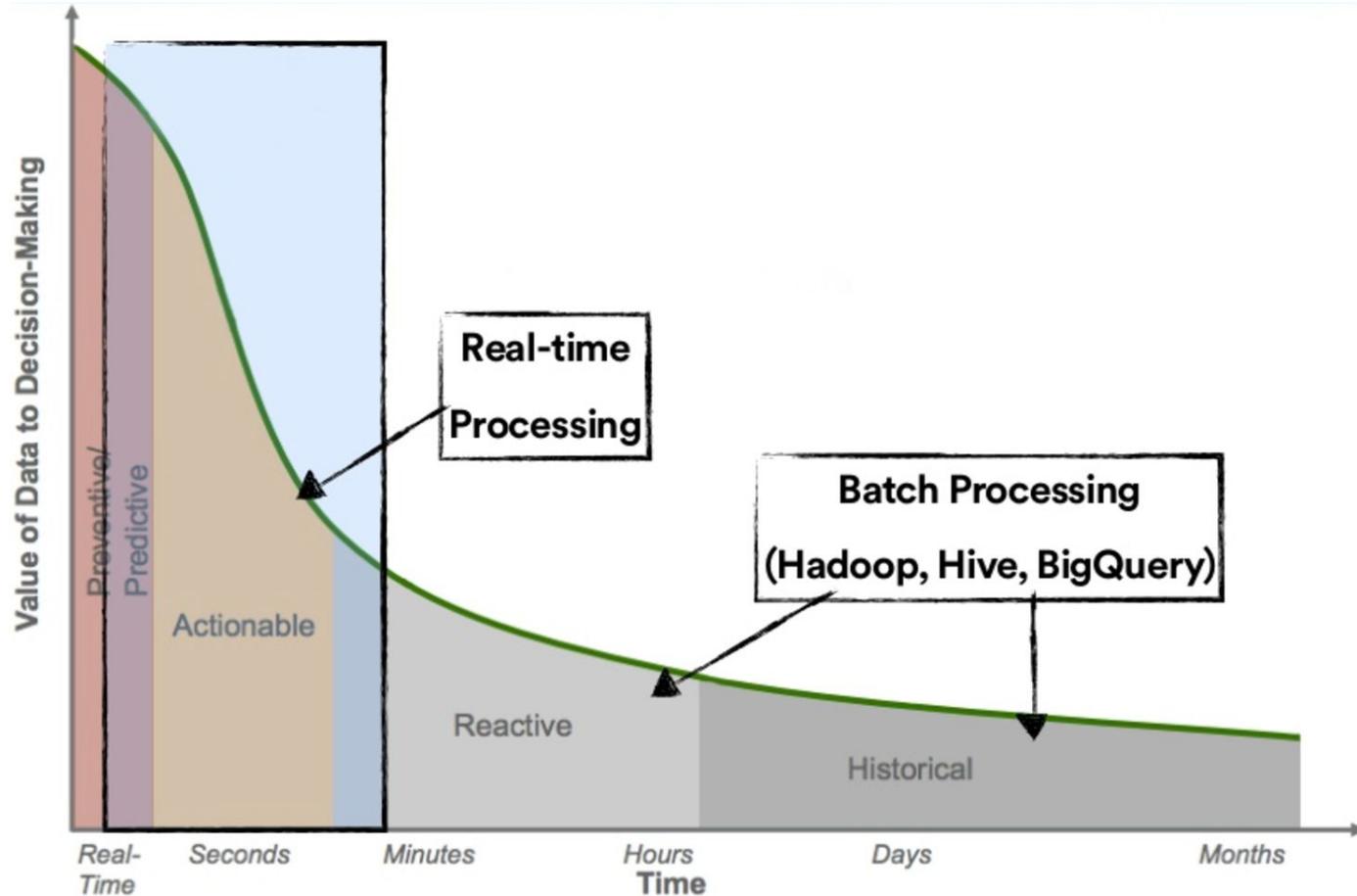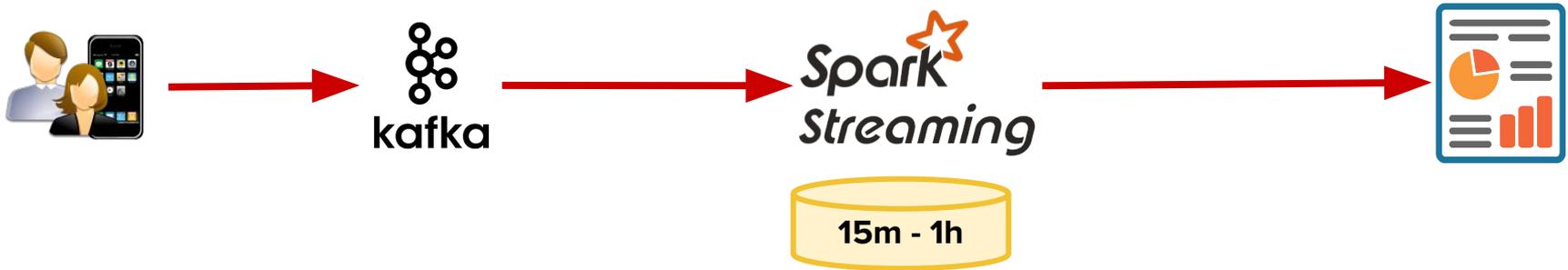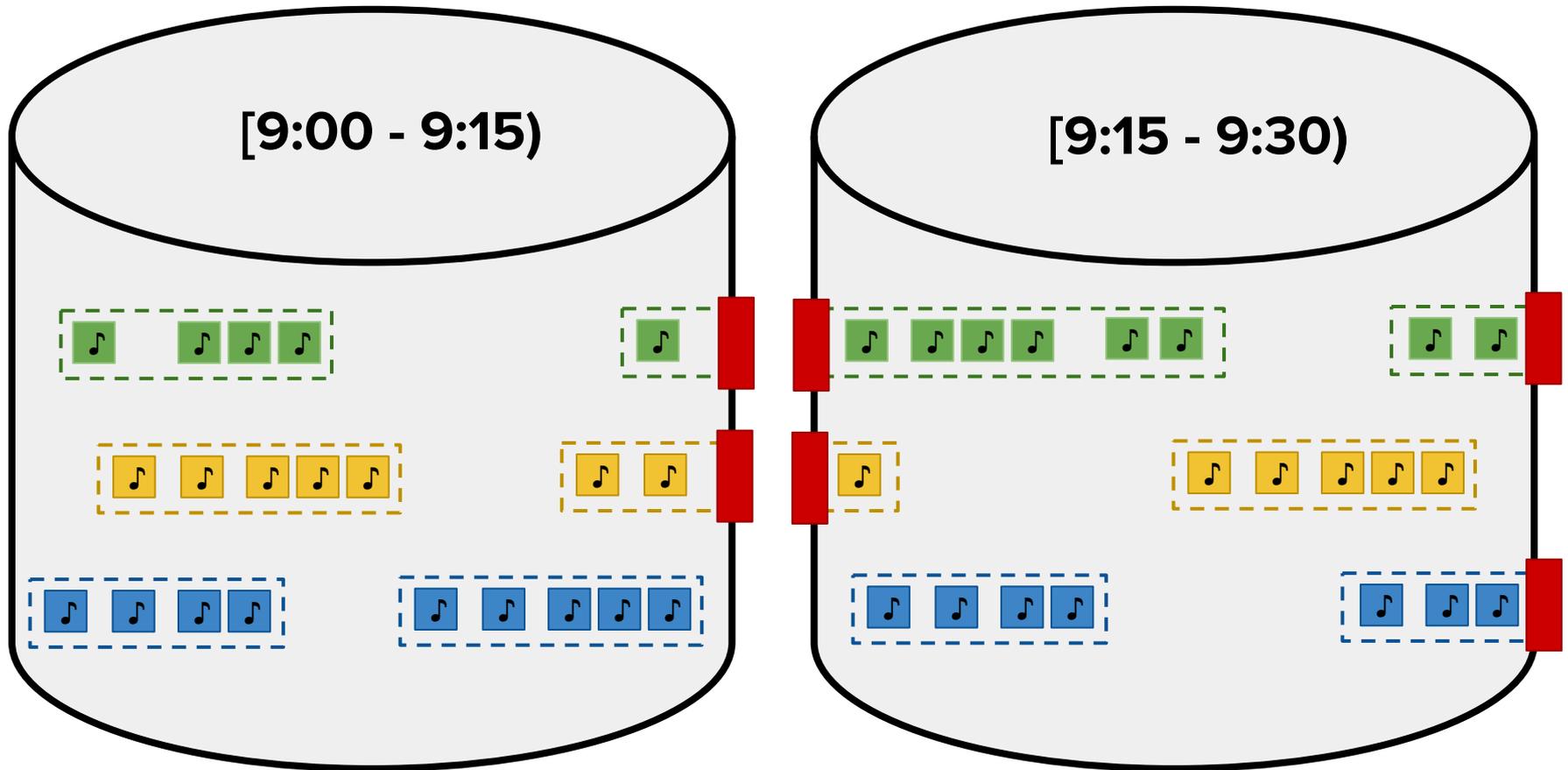
⬇ **The more error-prone solution**

# Long Waiting Time

# 2<sup>nd</sup> - Micro-Batch Architecture

kafka → Spark Streaming

15m - 1h

# No Built-In Session Windows



[9:00 - 9:15)

[9:15 - 9:30)

# No Built-In Session Windows

**[9:00 - 9:15**

Faster Stateful Stream Processing in Apache Spark Streaming

by Tathagata Das and Shixiong Zhu
Posted in **ENGINEERING BLOG** | February 1, 2016
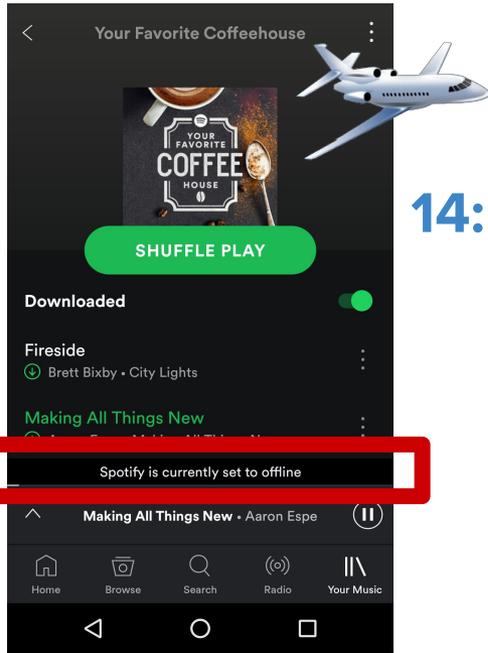
**Exploring Stateful Streaming with Apache Spark**

31 Jul 2016

**FIND OUT MORE!**

How-to: Do Near-Real Time Sessionization with Spark Streaming and Apache Hadoop

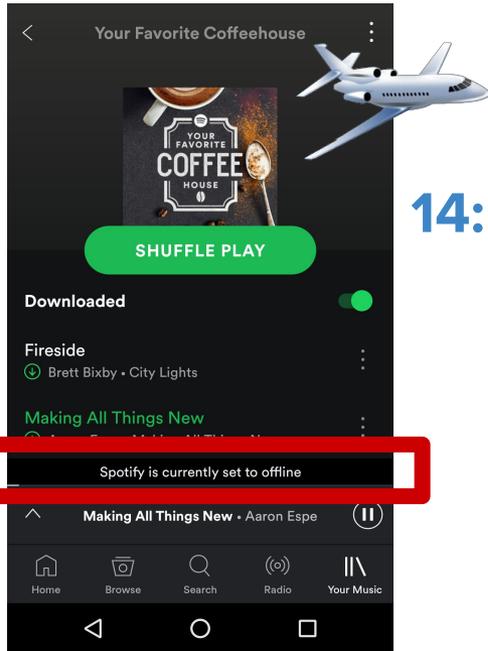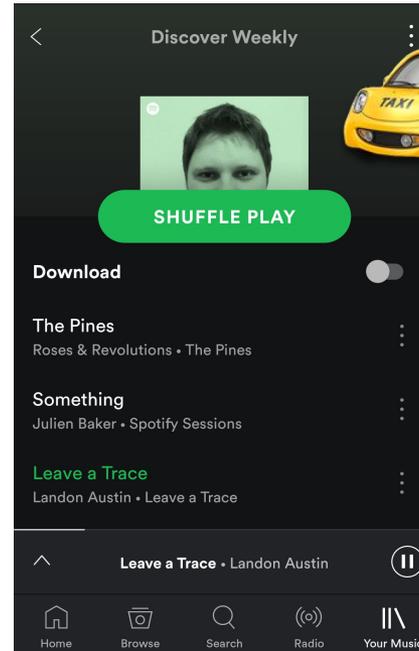November 3, 2014 | By Ted Malaska | 12 Comments

# Late Data ...



14:55 - 16:45
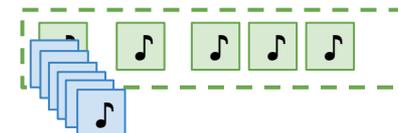
Event Time

Processing Time

# ... Included in Current Batch



14:55 - 16:45

16:55 - ...

Event Time

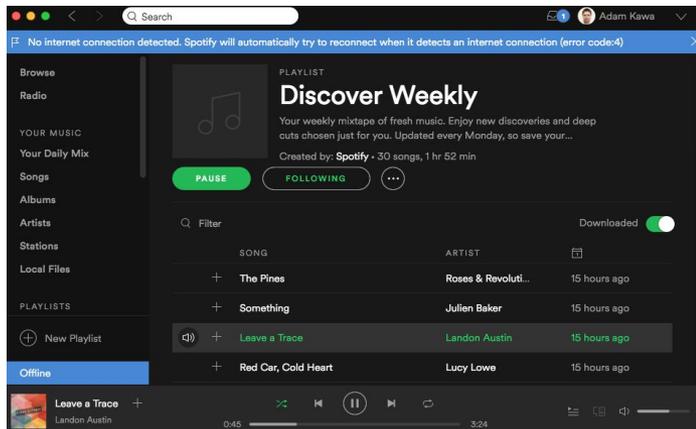Processing Time

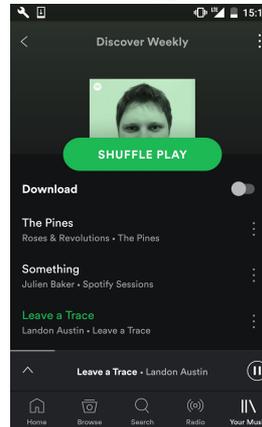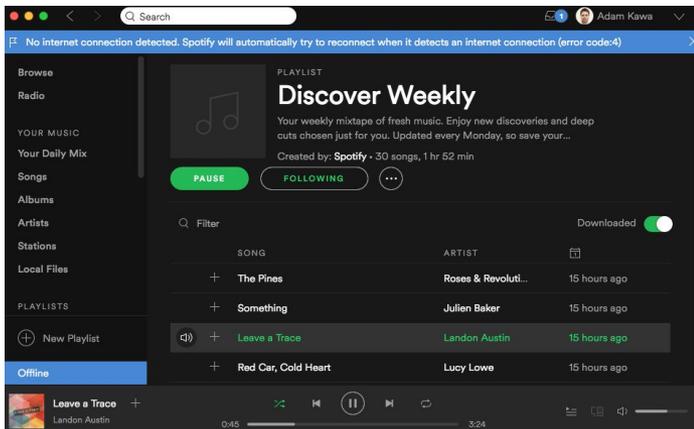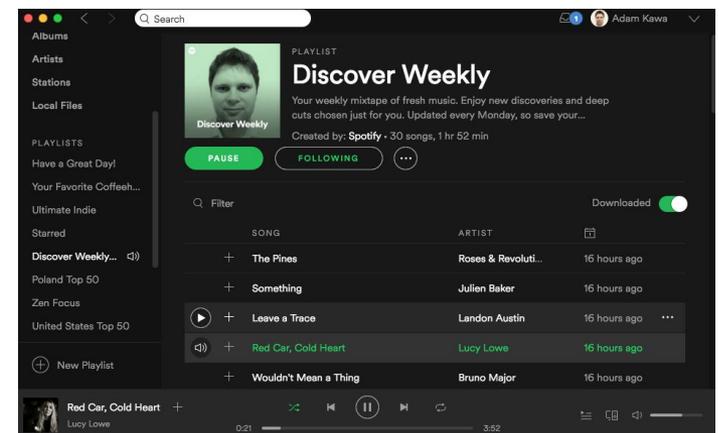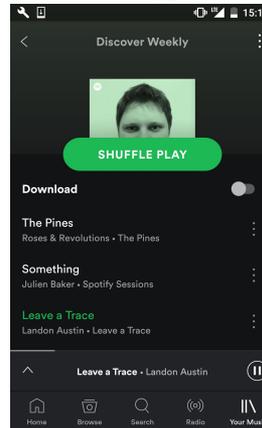# Out-Of-Order Data ...

**Event Time**

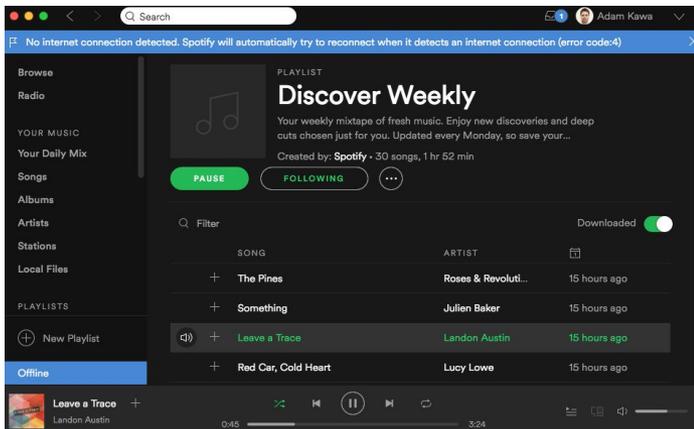**Processing Time**

# Out-Of-Order Data ...
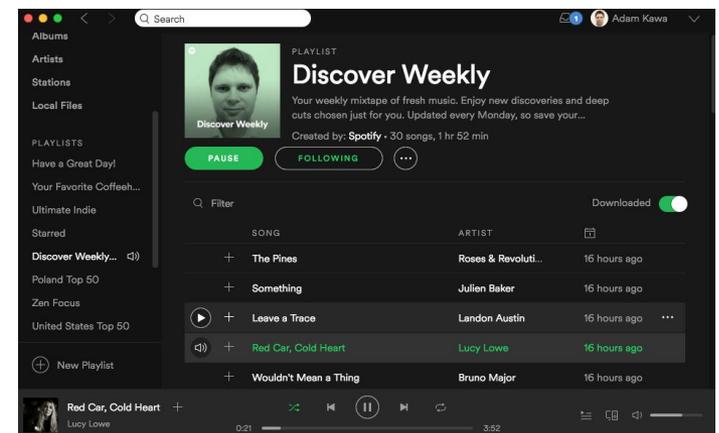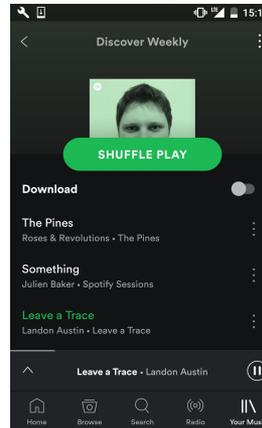


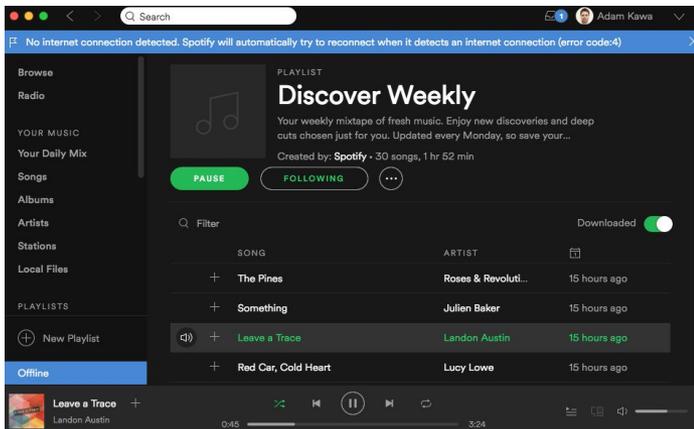Event Time

Processing Time

# Out-Of-Order Data ...



Event Time

Processing Time

# ... Breaks Correctness



Event Time

Processing Time

# Problems

**FILES,
    BATCHES,
    DATA LAKES**

# Solving Streaming Problem With Batch?

# 3<sup>rd</sup> - Streaming-First Architecture

# User Session Windows



Case A

Case B

Session gap
eg. 15
minutes

# Reading From Kafka

```
val sessionStream : DataStream[SessionStats] = sEnv
    .addSource(new KafkaConsumer(...))
```
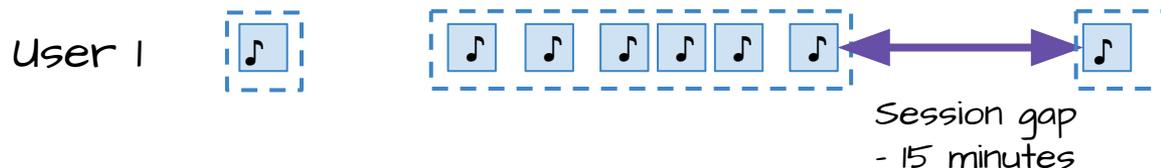
# Session Windows With Gap

```
val sessionStream : DataStream[SessionStats] = sEnv
    .addSource(new KafkaConsumer(...))
    .keyBy(_.userId)
```

User 1   ♪ ♪ ♪       ♪  ♪ ♪ ♪ ♪   ♪

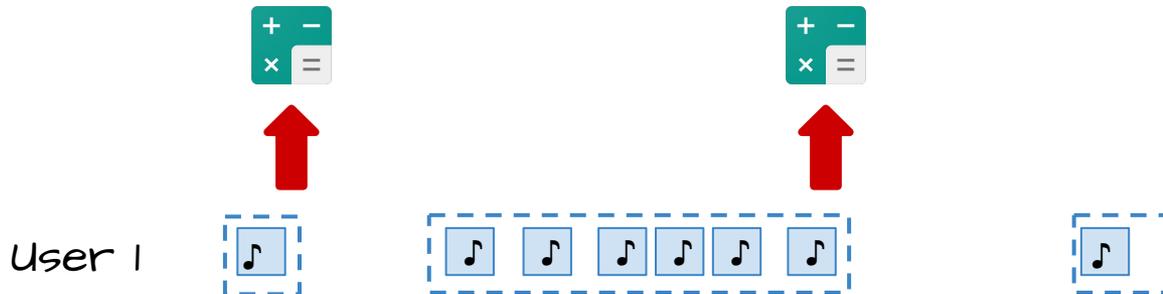User 2       ♪ ♪           ♪   ♪ ♪ ♪ ♪

# Session Windows With Gap

```
val sessionStream : DataStream[SessionStats] = sEnv
    .addSource(new KafkaConsumer(...))
    .keyBy(_.userId)
    .window(EventTimeSessionWindows.withGap(Time.minutes(15)))
```
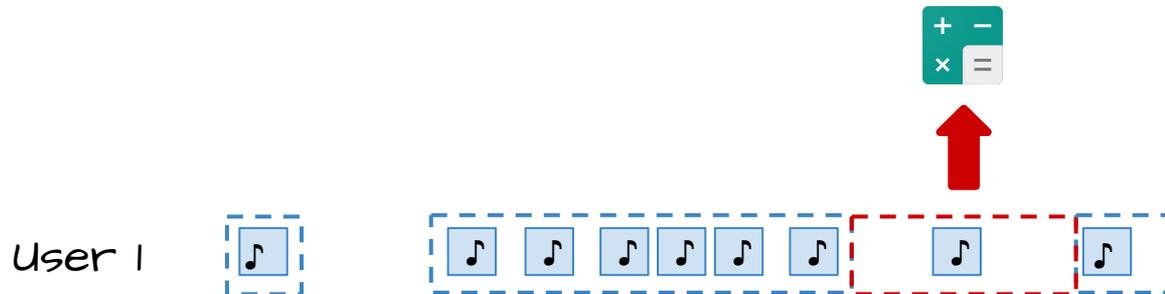
User 1

Session gap
- 15 minutes

# Analyzing User Session

```
val sessionStream : DataStream[SessionStats] = sEnv
    .addSource(new KafkaConsumer(...))
    .keyBy(_.userId)
    .window(EventTimeSessionWindows.withGap(Time.minutes(15)))
    .apply(new CountSessionStats())
```
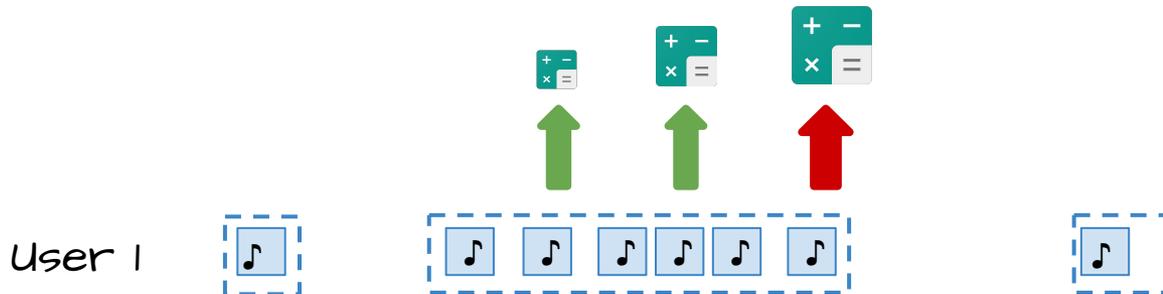
# Handling Late Events

```scala
val sessionStream : DataStream[SessionStats] = sEnv
    .addSource(new KafkaConsumer(...))
    .keyBy(_.userId)
    .window(EventTimeSessionWindows.withGap(Time.minutes(15)))
    .allowedLateness(Time.minutes(60))
    .apply(new CountSessionStats())
```

User 1

# Triggering Early Results

```
val sessionStream : DataStream[SessionStats] = sEnv

    .addSource(new KafkaConsumer(...))

    .keyBy(_.userId)

    .window(EventTimeSessionWindows.withGap(Time.minutes(15)))

    .trigger(EarlyTriggeringTrigger.every(Time.minutes(10)))

    .allowedLateness(Time.minutes(60))

    .apply(new CountSessionStats())
```

# Sessionization Example

```scala
val sessionStream : DataStream[SessionStats] = sEnv
    .addSource(new KafkaConsumer(...))
    .keyBy(_.userId)
    .window(EventTimeSessionWindows.withGap(Time.minutes(15)))
    .trigger(EarlyTriggeringTrigger.every(Time.minutes(10)))
    .allowedLateness(Time.minutes(60))
    .apply(new CountSessionStats())
```

FIND OUT MORE!

**Working example:**
**https://github.com/getindata/flink-use-case**

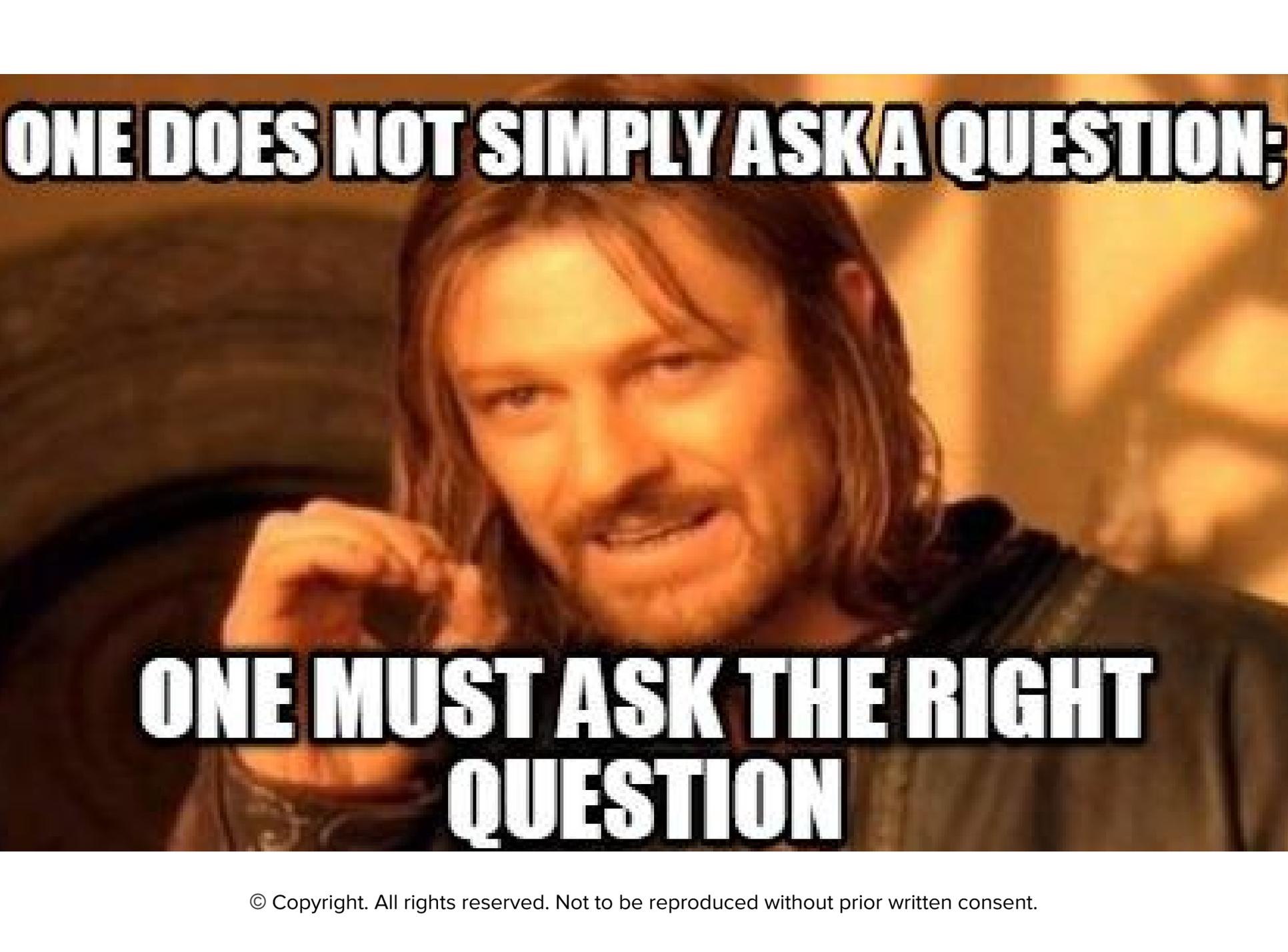# Modern Stream Processing Engines

- **Rich stream processing semantic**
  - Built-in support for event-time windows
  - Accurate results for late / out-of-order events and replays
  - Early triggers
- **Low latency and high-throughput**
- **Exactly-once stateful processing**

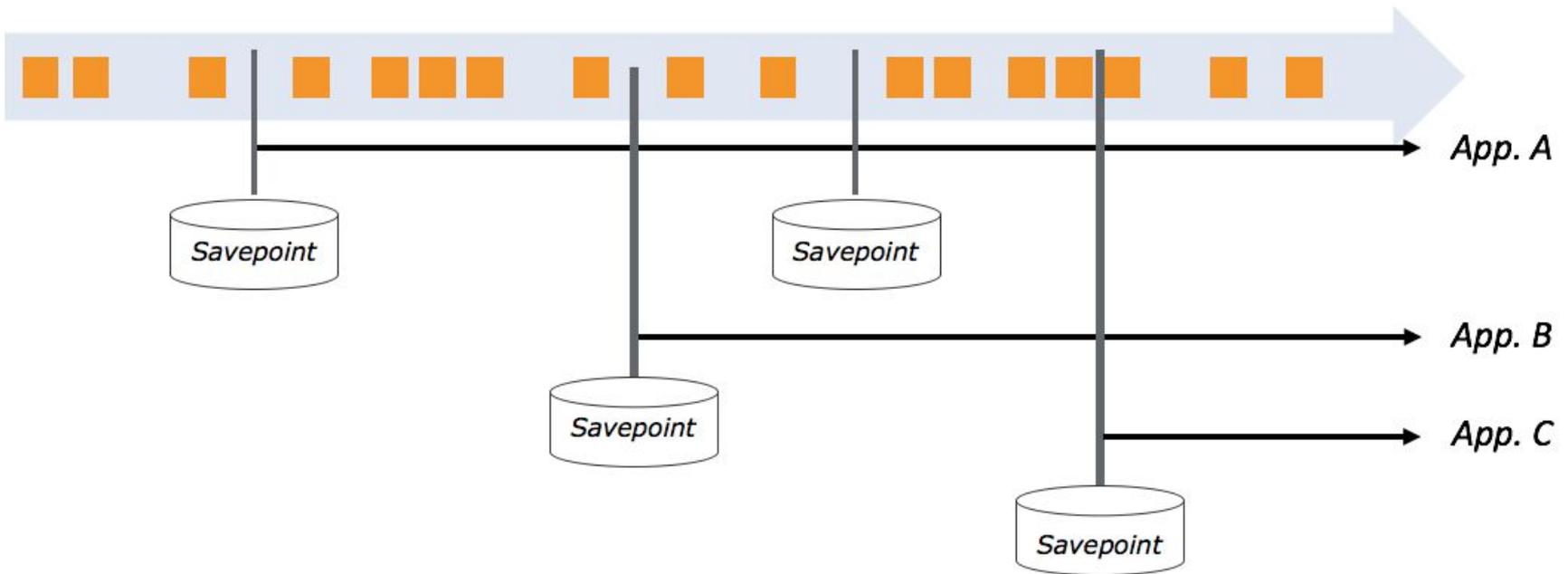# Modern Stream Processing Engines

- **Rich stream processing semantic**
  - Built-in support for event-time windows
  - Accurate results for late / out-of-order events and replays
  - Early triggers
- **Low latency and high-throughput**
- **Exactly-once stateful processing**

**FIND OUT MORE!**

**User survey:**

**http://data-artisans.com/flink-user-survey-2016-part-1**
**http://data-artisans.com/flink-user-survey-2016-part-2**

ONE DOES NOT SIMPLY ASK A QUESTION; ONE MUST ASK THE RIGHT QUESTION

# How can I **<span style="color:red">reprocess</span>** data?

# Reprocessing Events In Flink

1. **Take periodic snapshots of a job**
   - It stores Kafka offsets, on-flight sessions, application state
2. **Restart a job from a savepoint rather than from a beginning**

# What if data is **no longer in Kafka**?

# Consuming Data From HDFS

- **Run your streaming code on HDFS (bounded data)**
  - You need to read data in event-time based order
  - Implement mechanism of proper watermark generation

# What are usual <span style="color:red">stream</span> processing applications?

# Stream Analytics

# Stream 24/7 Applications

# What can happen
# if you <span style="color:red">stick to batch</span> ?

# Real-Time Personalization



Image source:

https://www.slideshare.net/g33ktalk/spotifys-ad-targeting-infrastructure-achieving-realtime-personalization-for-2 4-million-users
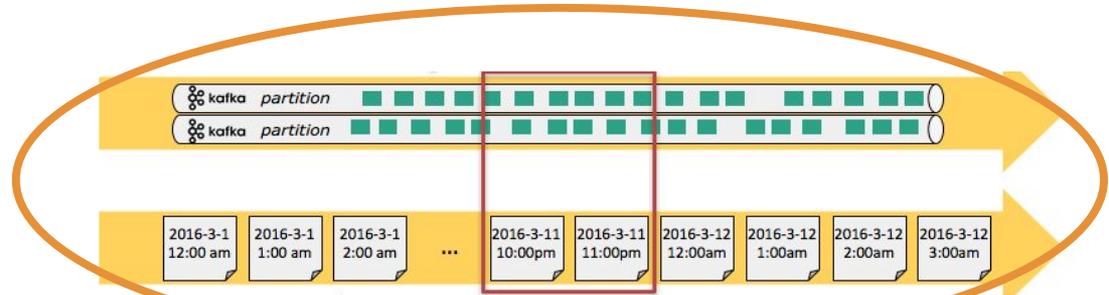
# When is batch processing good?

# Batch Processing Use-Cases

- **Ad-hoc analytics and data exploration**
  - Notebooks, Spark/**Flink**/Hive, Parquet
  - Complete data sets
- **Technical advantages**
  - A large swaths of historical data in HDFS
  - High-level libraries in mature batch technologies

# Batch or Stream?

- **Stream is often a natural representation of data**
- **Stream processing is not only about low latency**
  - Correctness, expressive API, simplicity
- **Stream processing is a great fit for ETL, KPIs, reports**

# Take-Away Message

- Stream is often a natural representation of data **(when data arrives continuously)**
- Stream processing is not only about low latency
  - Correctness, expressive API, simplicity
- Stream processing is a great fit for ETL, KPIs, reports

**don't solve streaming problem with batch jobs**

# Thanks Vilnius !

# Q&A



big data. experience. passion.