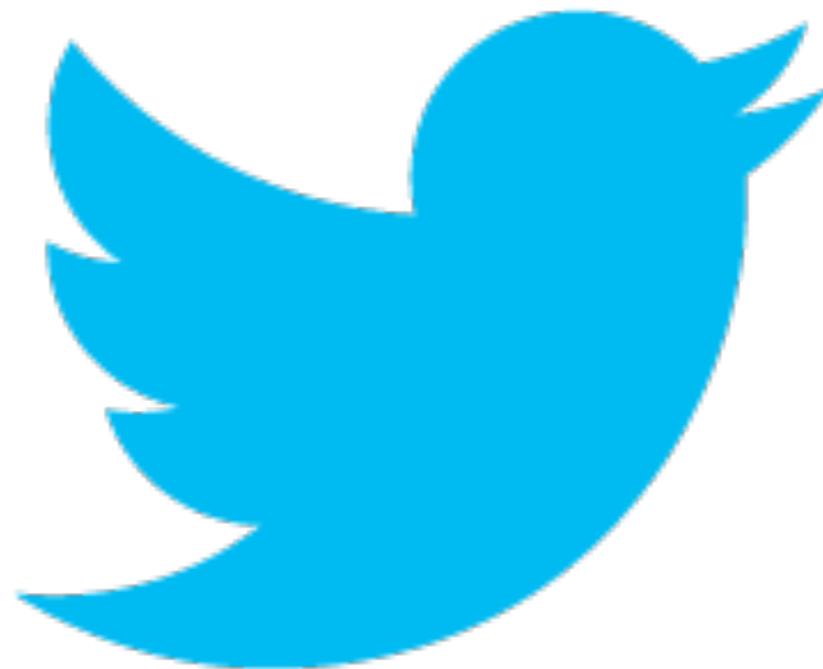# The Road to a Complete Tweet Index

Michael Busch
@michibusch
michael@twitter.com
buschmi@apache.org

# Introduction

**More than 2 billion search queries per day.**

# Introduction

**500 million** tweets are sent per day.

# Introduction

**Hundreds of billions of** tweets have been sent since company founding in 2006.

# 2010

# 2010

Realtime Search powered by Summize technology

# 2010

Realtime Search powered by Summize technology

# 2011

# 2011

Twitter launches first Lucene-based search engine: **Earlybird.**
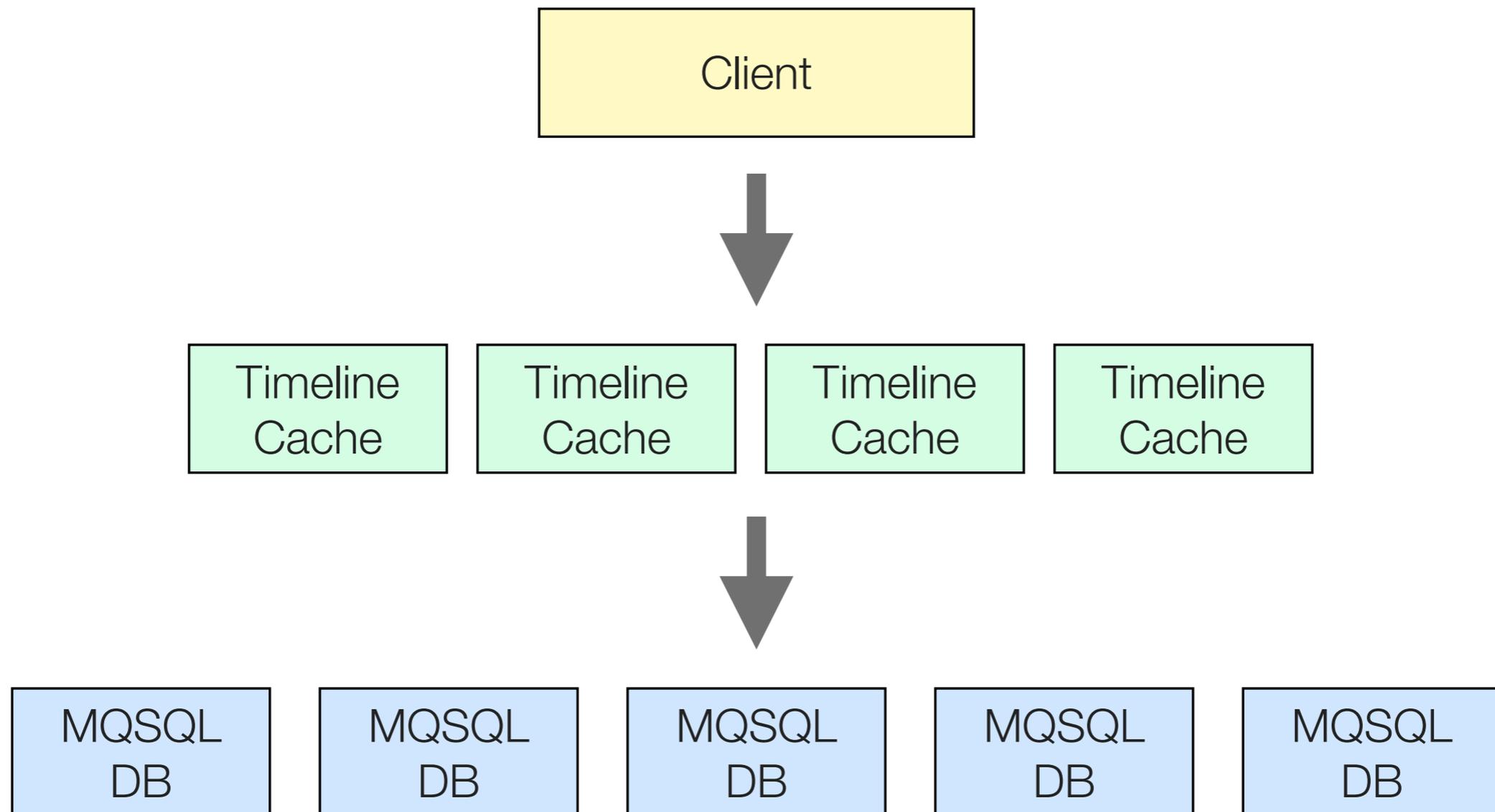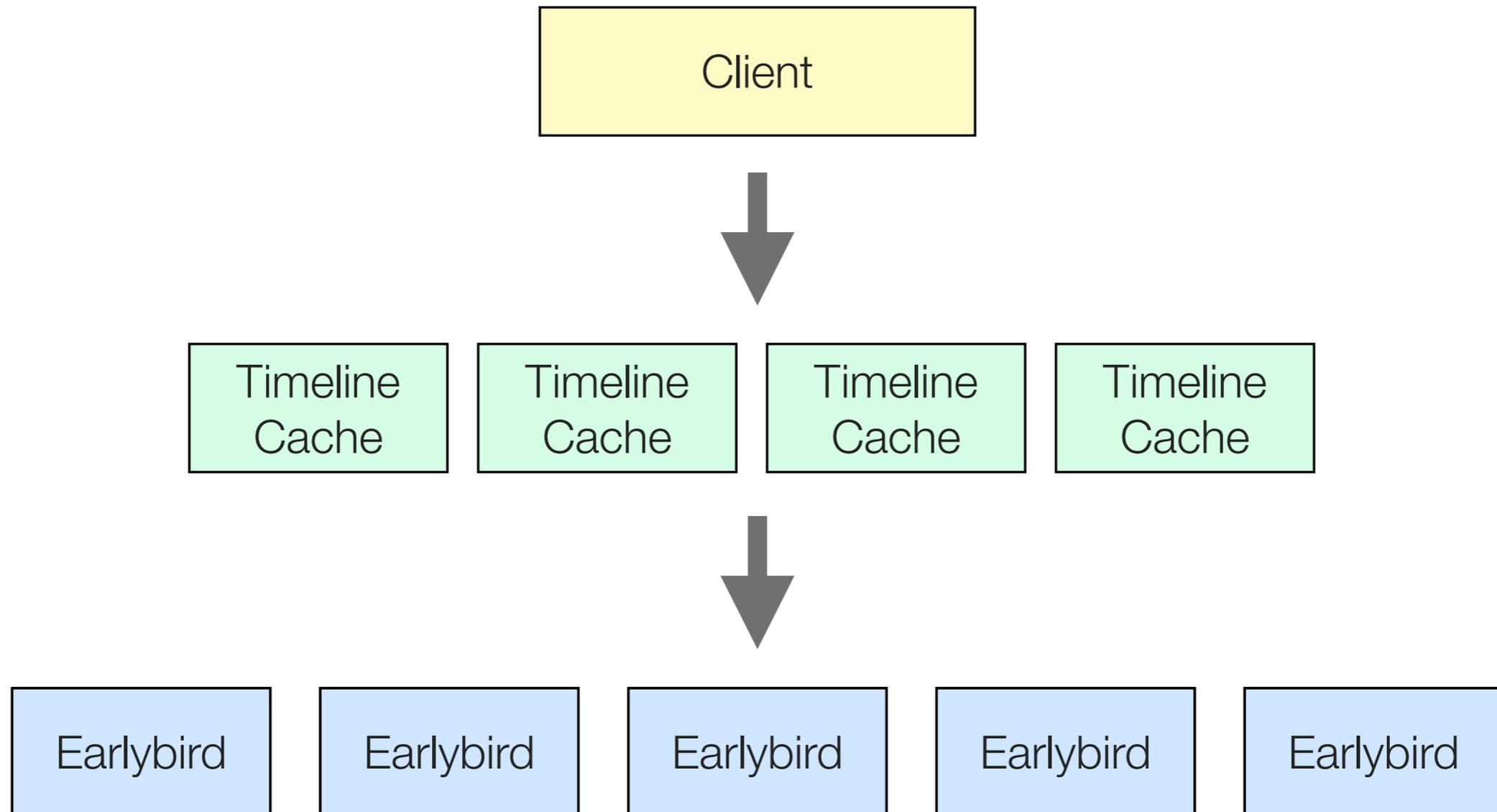
# 2011

Realtime Search powered by Summize technology

# 2011

Realtime Search powered by Earlybird
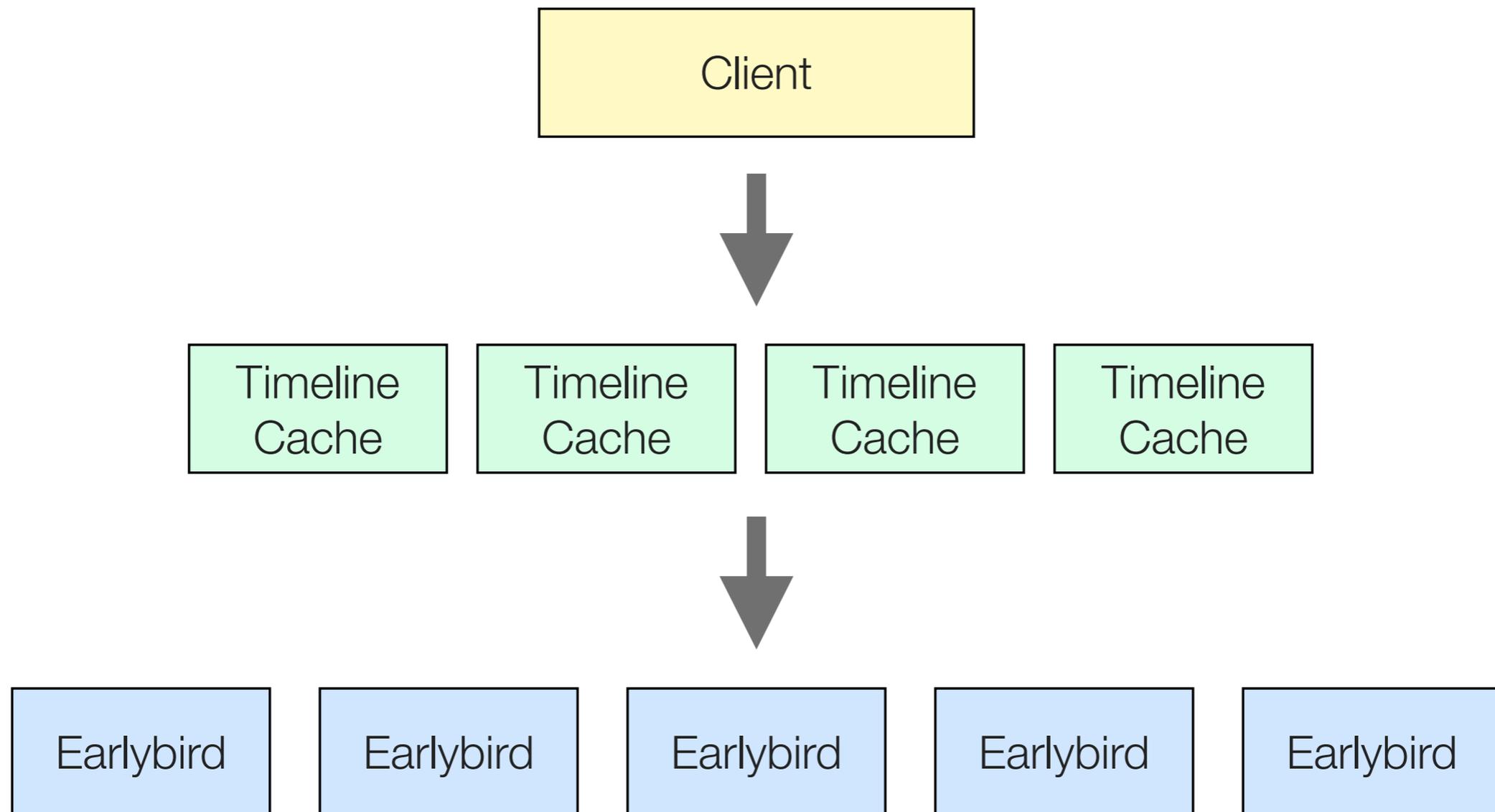
# Earlybird

- In-memory index containing several days of most recent tweets

- Highly optimized for realtime search

- Limited to short documents (max. 255 tokens)

- Novel concurrency and memory models

- Concurrently writing and searching an index segment
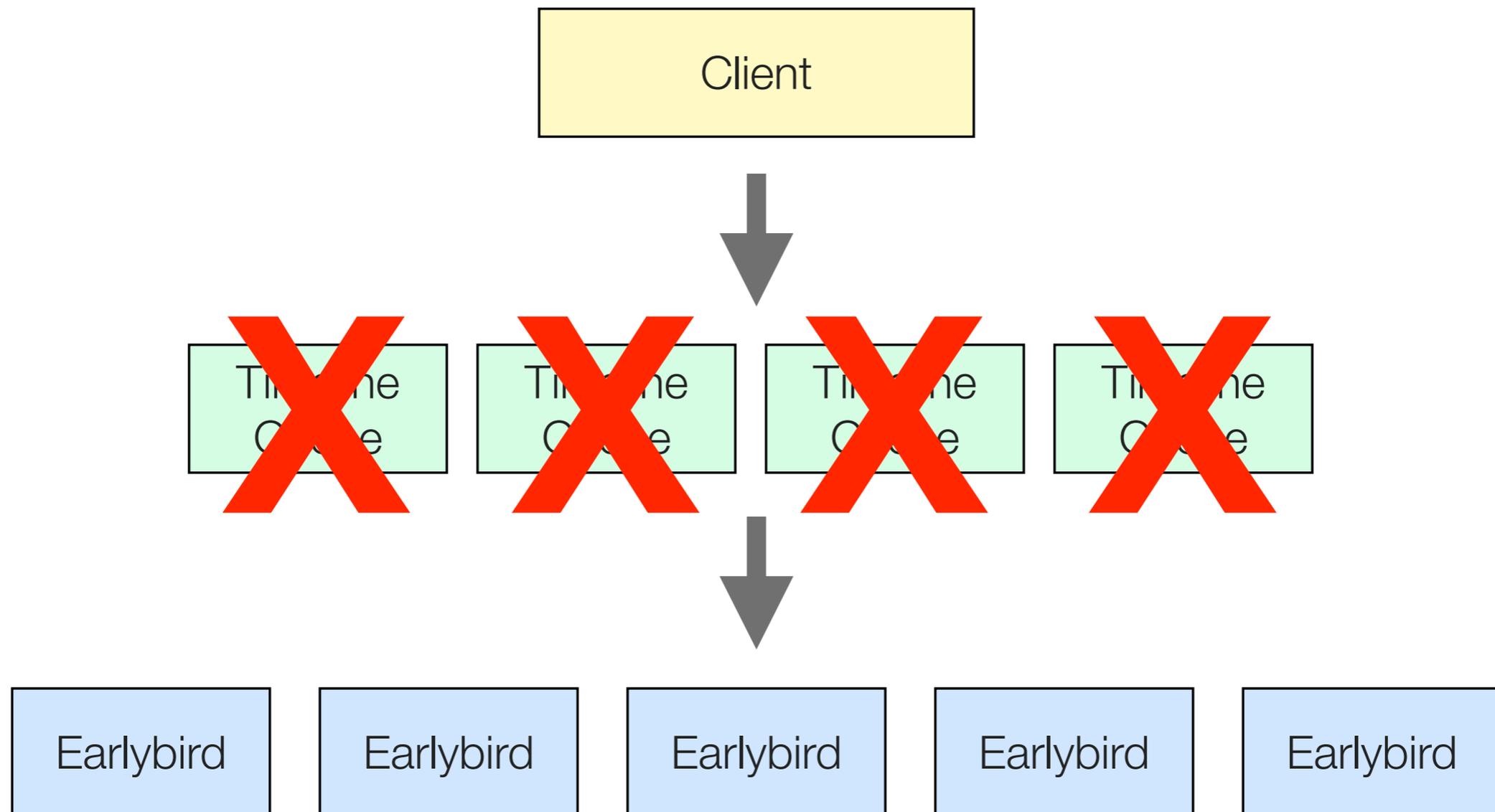
# 2011

Realtime Search powered by Earlybird

# 2011

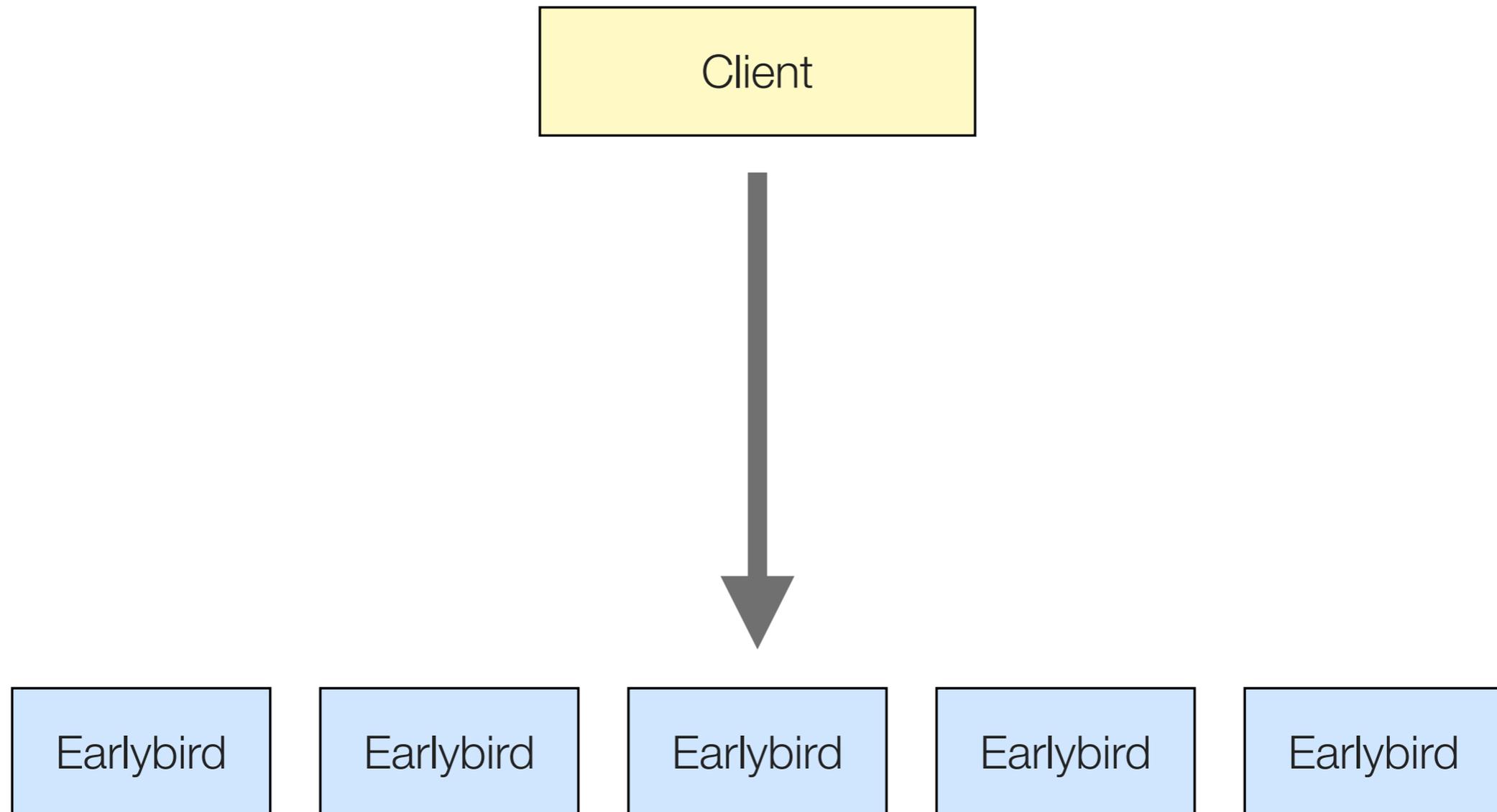Realtime Search powered by Earlybird

# 2011

Realtime Search powered by Earlybird

# Earlybird

Realtime indexing pipeline



| | raw tweets<br>(JSON) | Analyzer/<br>Partitioner | analyzed tweets<br>(Thrift) | Earlybird<br>Indexes |

# Earlybird



*m **hash partitions**; hash function is simply document uid % m*

# Earlybird



m **hash partitions**; hash function is simply document uid % m

# Earlybird

# Earlybird

Fixed segment sizes



Indexing direction

Search direction

Complete segment (8M tweets)

Current incomplete segment

# Earlybird

Fixed segment sizes



Indexing direction

Search direction

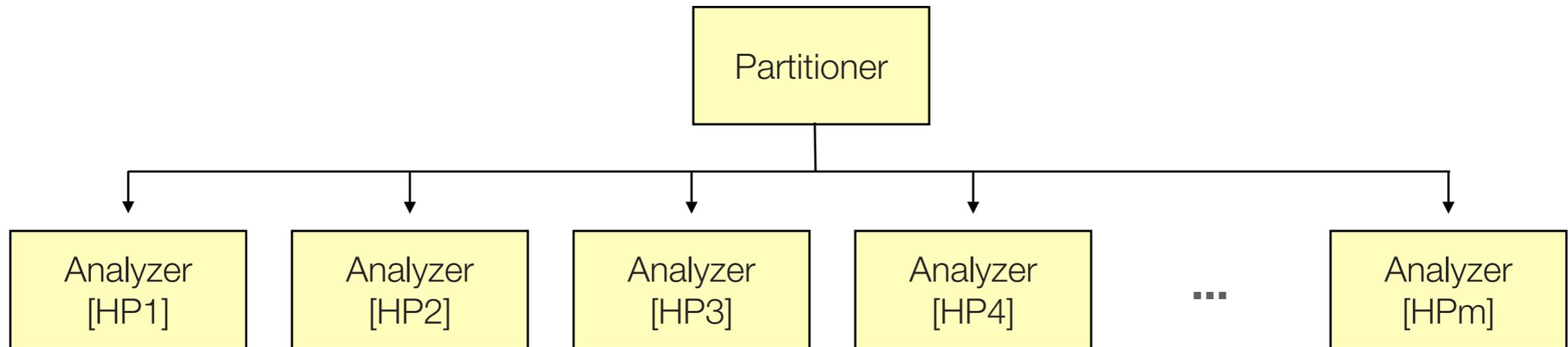**Current segment** *can be* **searched** *while its* **being built***; no segment flushing/merging necessary; lock-free concurrency model*

Complete segment (8M tweets)

Current incomplete segment

# Earlybird

Fixed segment sizes



Indexing direction

Search direction

**Sliding window**: *number of segments per Earlybird index is constant;* **oldest segment** *is deleted when new segment is started*

Complete segment (8M tweets)

Current incomplete segment

# Earlybird

**Why use a fixed segment size?**

• Earlybird does not need to flush a segment to make it searchable

• No segment merges necessary - consistently high indexing throughput

• Predictable indexing and search performance

• Fixed segment sizes keep replicated Earlybird indexes in sync

# Earlybird



the **replicas** share complete segments via **HDFS**

# Earlybird

# Earlybird



Partitioner

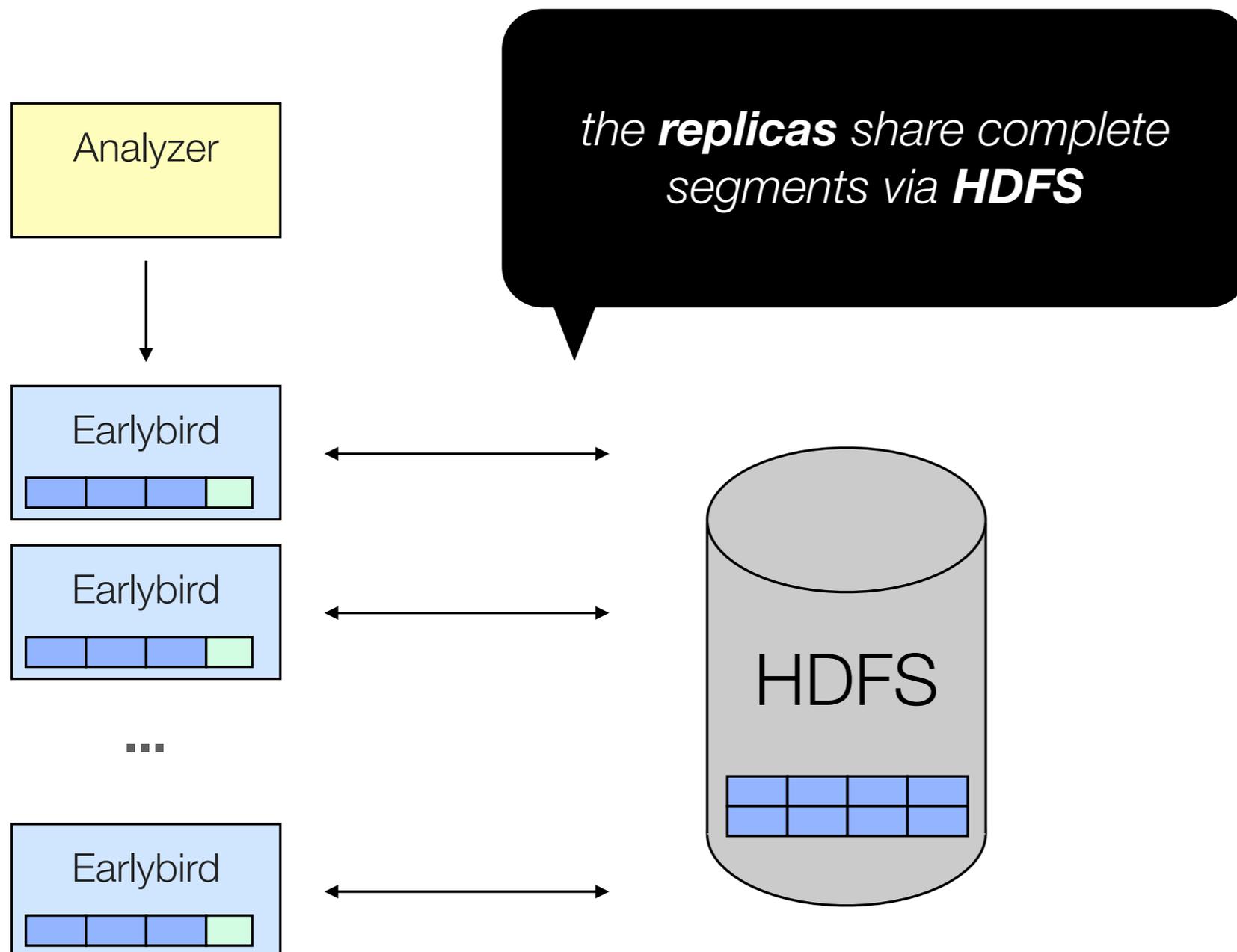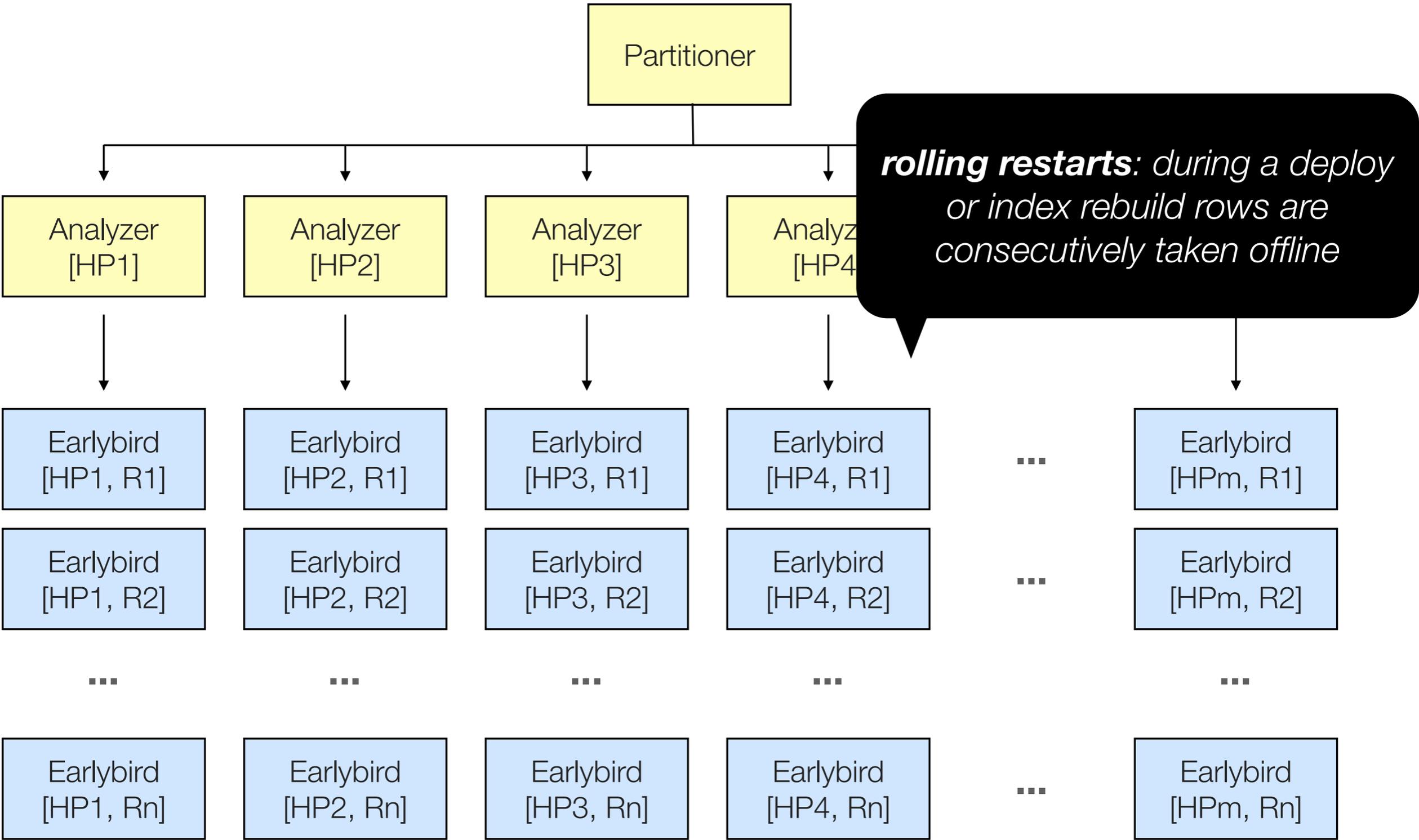Analyzer [HP1]   Analyzer [HP2]   Analyzer [HP3]   Analyz [HP4]

*the **first** replica in each hash partition **builds** complete segments and copies them to **HDFS***

Earlybird [HP1, R1]   Earlybird [HP2, R1]   Earlybird [HP3, R1]   Earlybird [HP4, R1]   ...   Earlybird [HPm, R1]

Earlybird [HP1, R2]   Earlybird [HP2, R2]   Earlybird [HP3, R2]   Earlybird [HP4, R2]   ...   Earlybird [HPm, R2]

...   ...   ...   ...   ...

Earlybird [HP1, Rn]   Earlybird [HP2, Rn]   Earlybird [HP3, Rn]   Earlybird [HP4, Rn]   ...   Earlybird [HPm, Rn]

# Earlybird



Partitioner

Analyzer [HP1]   Analyzer [HP2]   Analyzer [HP3]   Analyzer [HP4]   ...   Analyzer [HPm]

Earlybird [HP1, R1]   Earlybird [HP2, R1]   Earlybird [HP3, R1]   Earlybird [HP4, R1]   Earlybird [HPm, R1]

Earlybird [HP1, R2]   Earlybird [HP2, R2]   Earlybird [HP3, R2]   Earlybird [HP4, R2]   ...   Earlybird [HPm, R2]

...   ...   ...   ...   ...

Earlybird [HP1, Rn]   Earlybird [HP2, Rn]   Earlybird [HP3, Rn]   Earlybird [HP4, Rn]   ...   Earlybird [HPm, Rn]

*subsequent* replicas restart faster by *loading* complete segments from **HDFS**

# Earlybird

# Earlybird

**Why copy segments to HDFS?**

- Rolling restarts are performed to deploy to Earlybird clusters

- A full cluster restart takes a multiple of the time it takes to restart a single hash partition

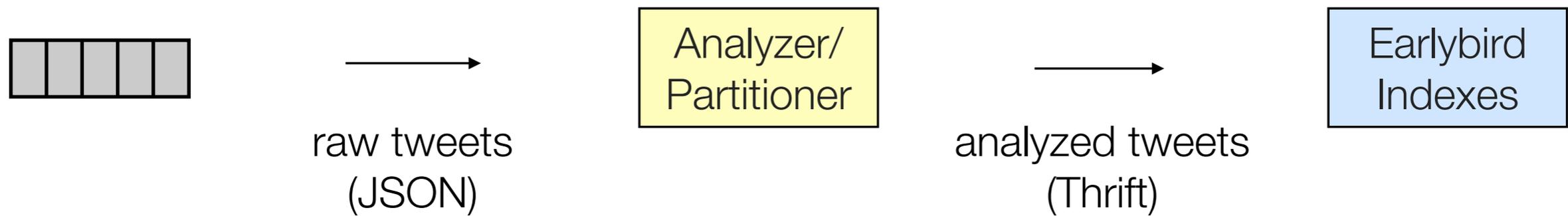- New machines can quickly bootstrap from HDFS

# 2012

# 2012

Twitter launches **in-memory historical index.**

# In-Memory historical index

- In-memory index containing approx. 2 billion top Tweets of **all time**

- Best Tweets per language

- "Best" determined by relevance function

- Inverted index format identical to realtime Earlybird
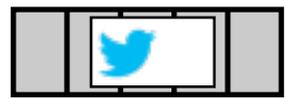
- New offline ingestion pipeline

# Ingestion pipeline

Online pipeline

# Ingestion pipeline

Online pipeline



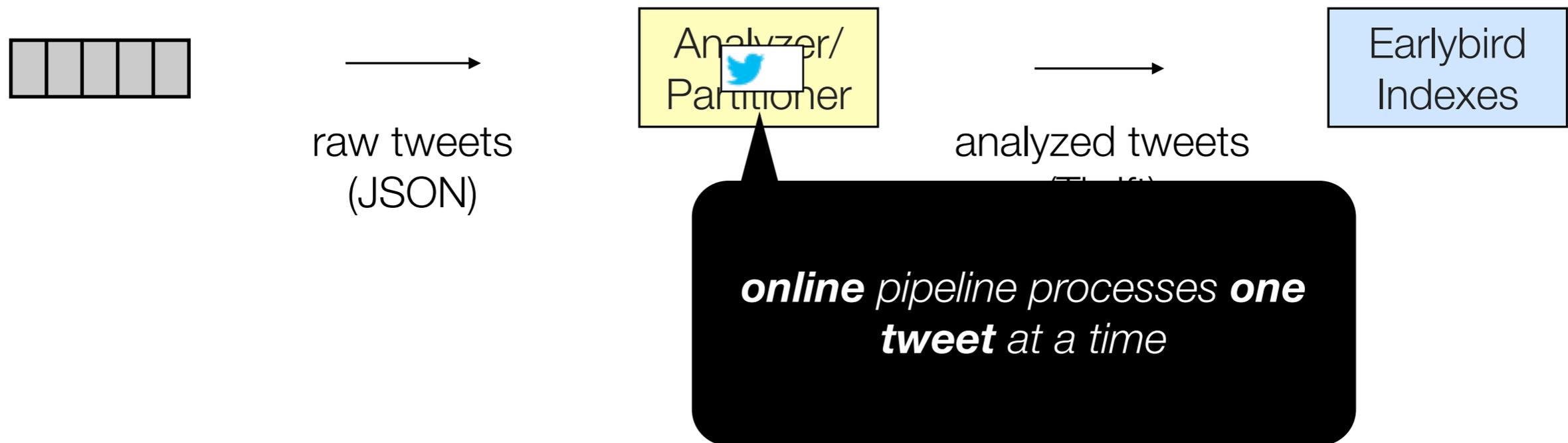raw tweets
(JSON)

Analyzer/
Partitioner

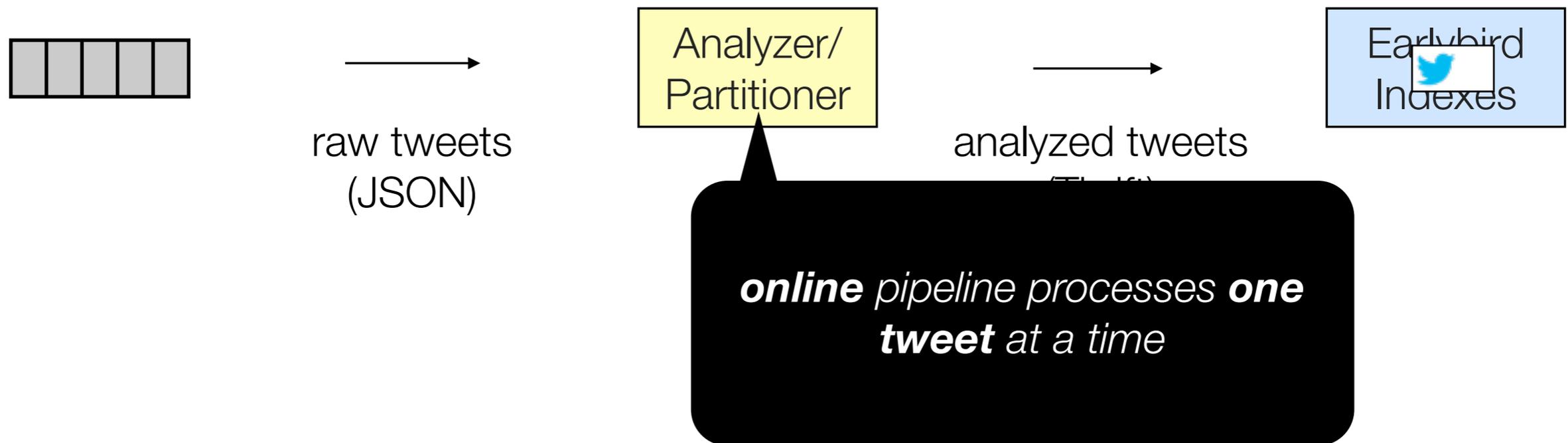analyzed tweets
(Thrift)

Earlybird
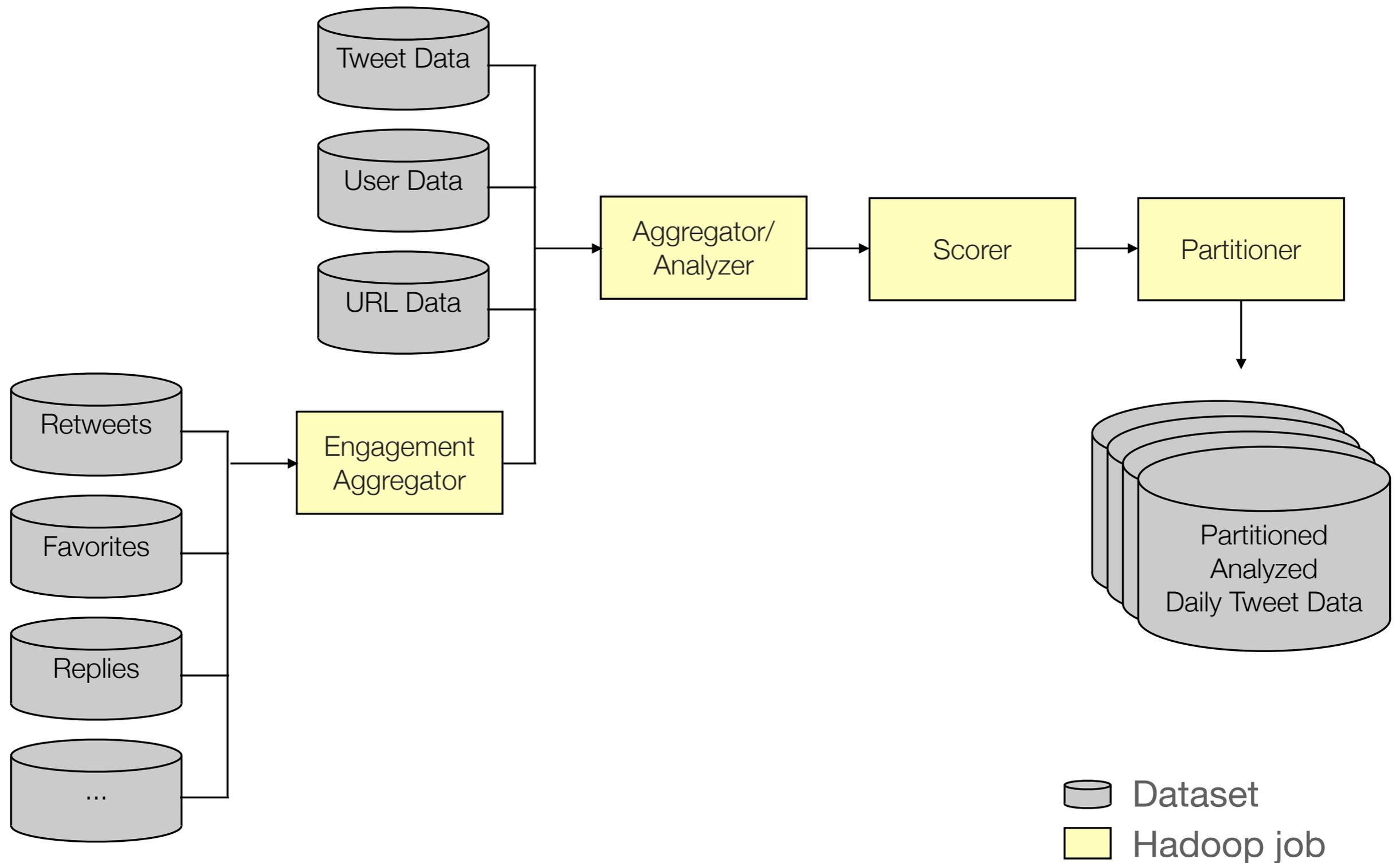Indexes

# Ingestion pipeline

Online pipeline



raw tweets
(JSON)

Analyzer/
Partitioner

analyzed tweets

Earlybird
Indexes

**online** *pipeline processes* **one tweet** *at a time*

# Ingestion pipeline

Online pipeline



raw tweets
(JSON)

Analyzer/
Partitioner

analyzed tweets

Earlybird
Indexes

**online** *pipeline processes* **one**
**tweet** *at a time*

# Offline Ingestion pipeline

# Offline Ingestion pipeline



Tweet Data

Scorer → Partitioner

**offline** *pipeline designed to process a* **daily batch** *of data at a time*

Retweets

Favorites

Replies

...

Engagement Aggregator

Partitioned Analyzed Daily Tweet Data
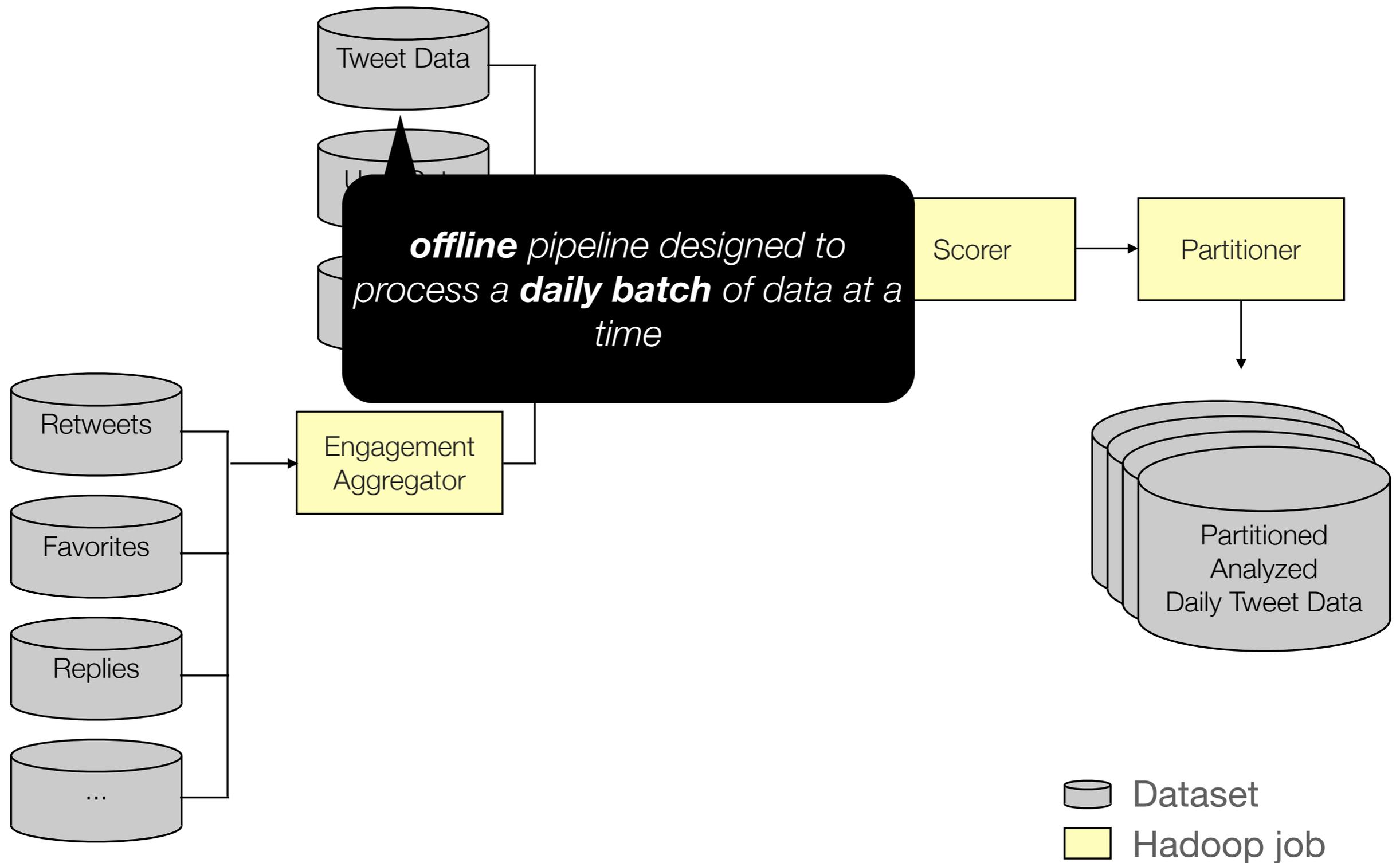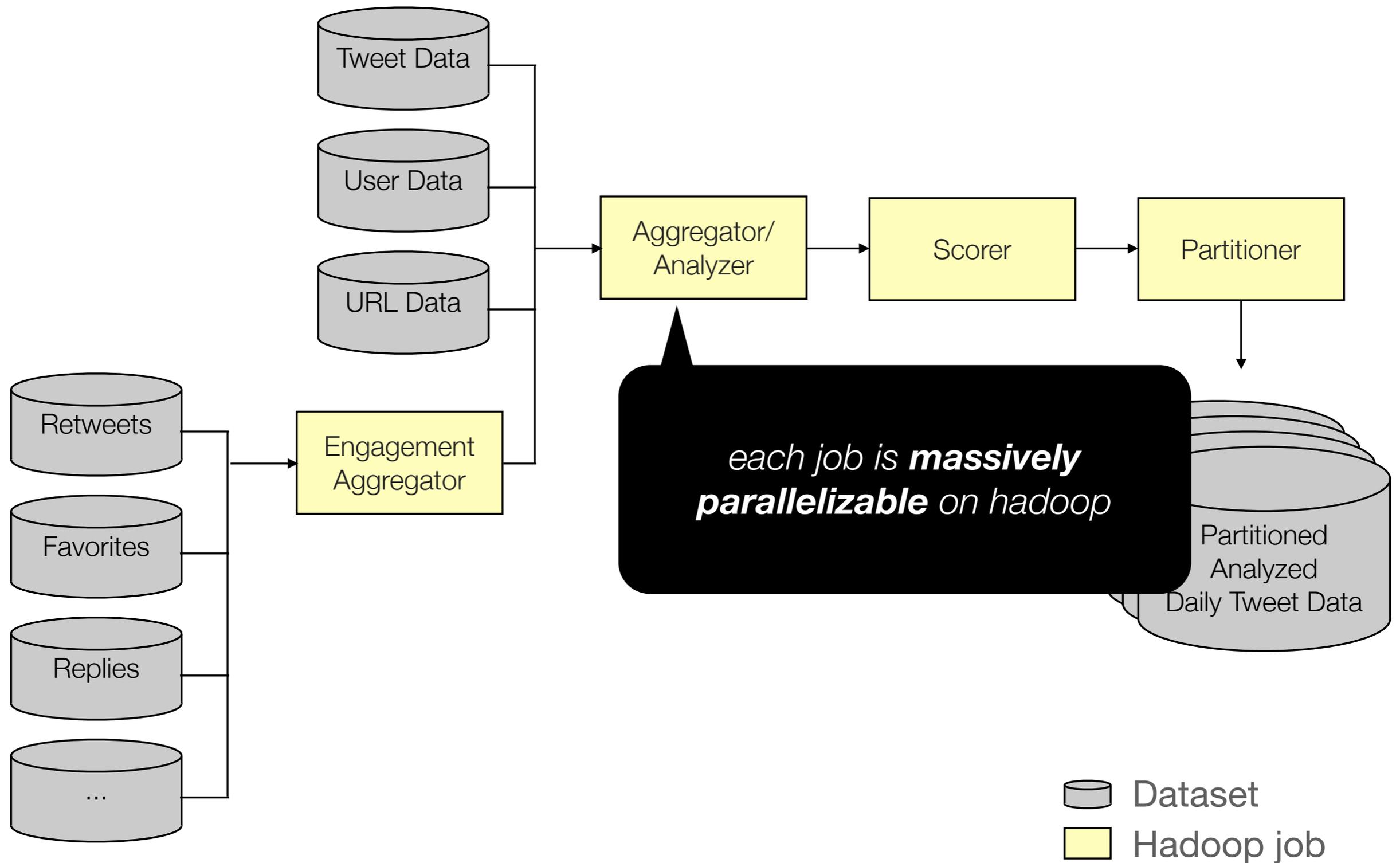
Dataset

Hadoop job
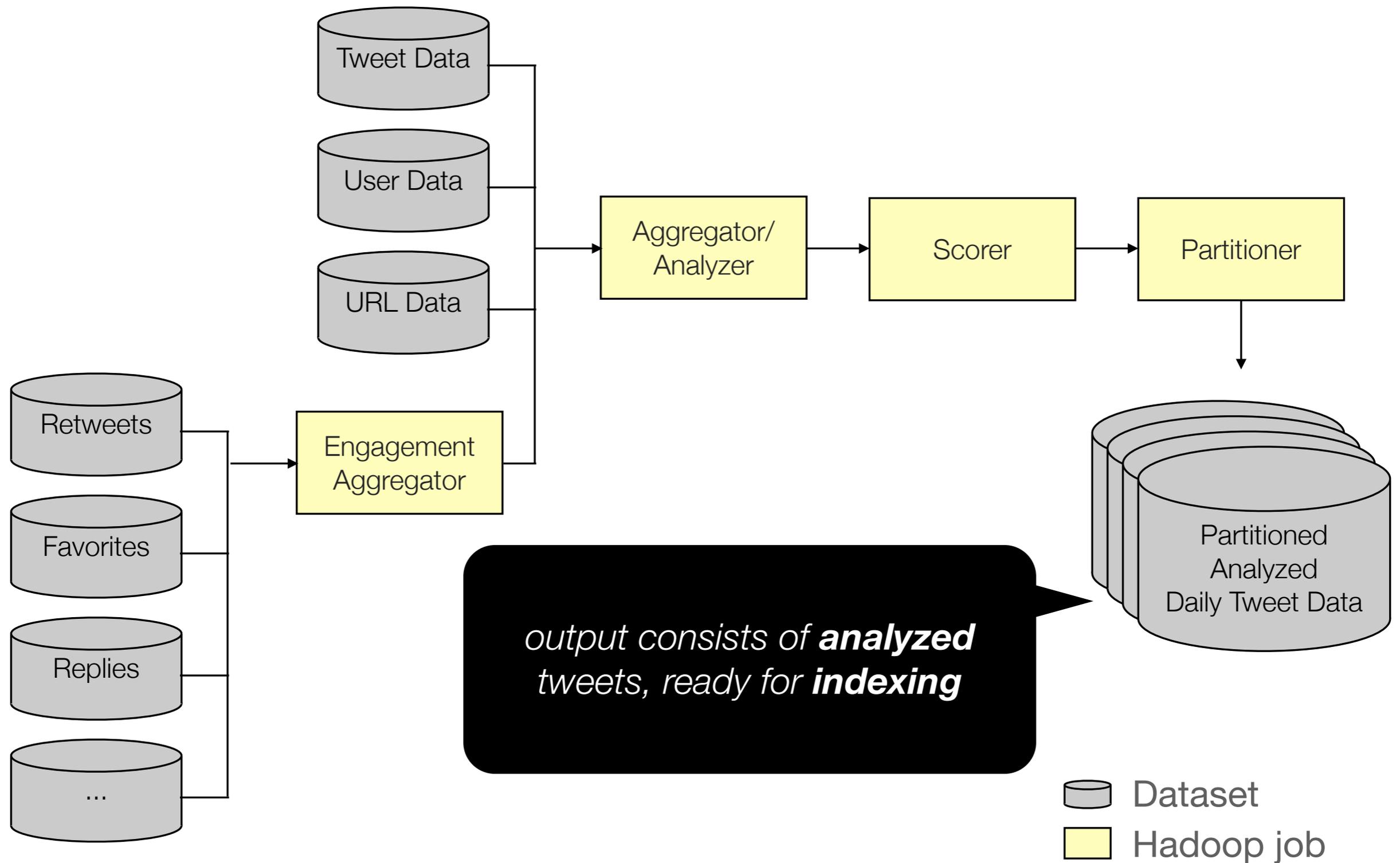
# Offline Ingestion pipeline
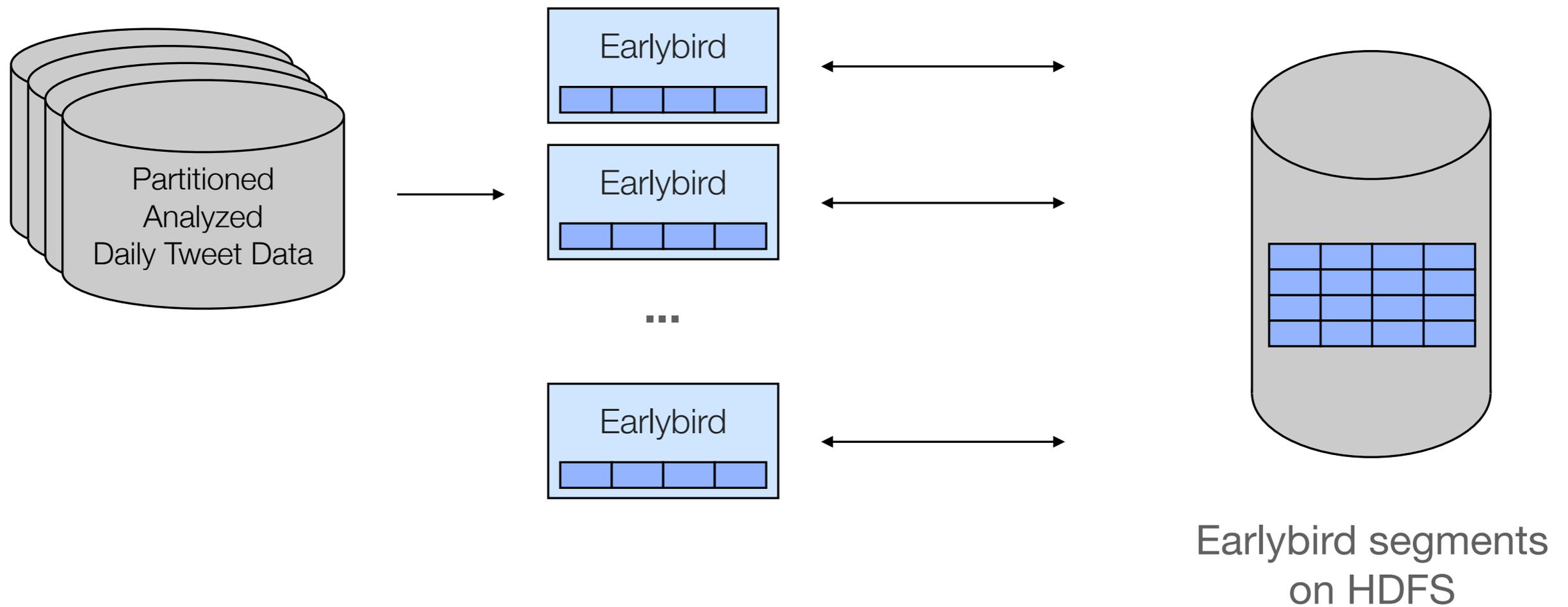
# Offline Ingestion pipeline
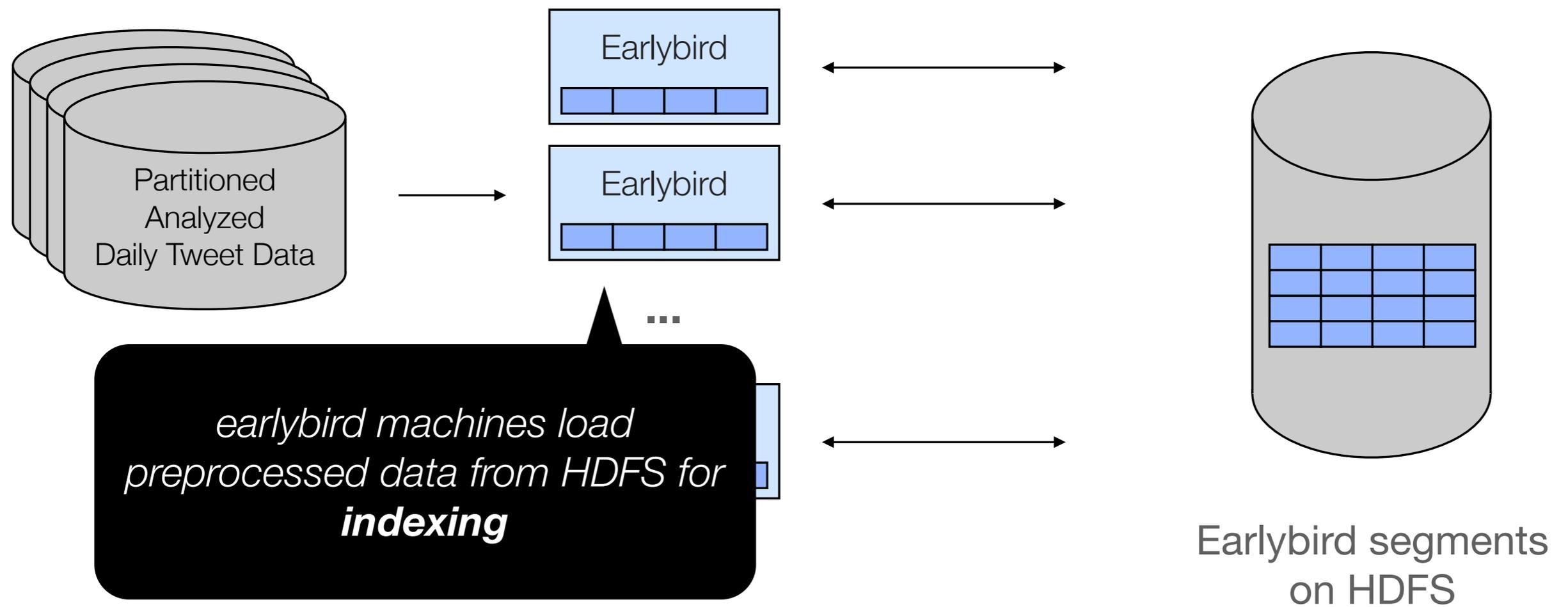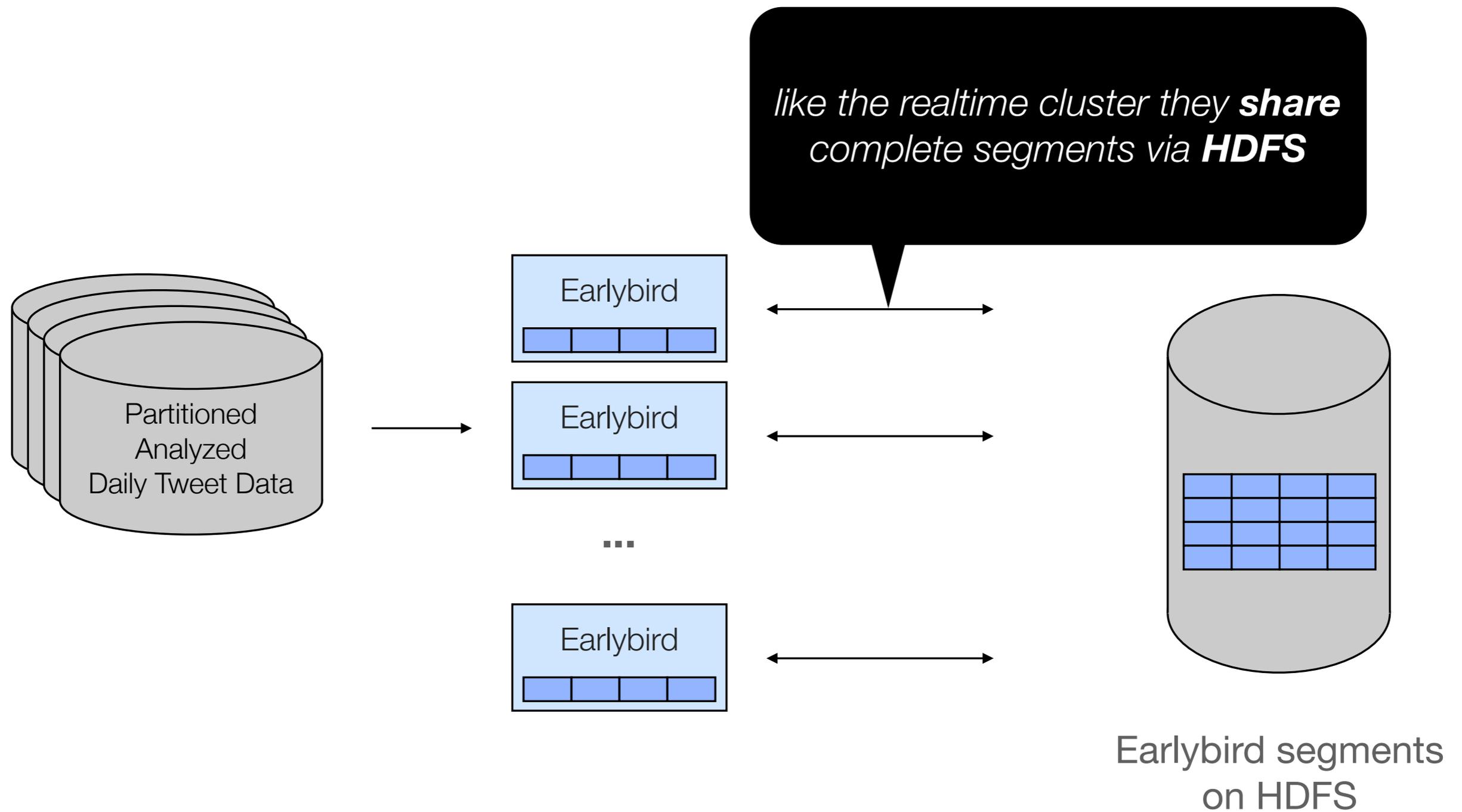
# Offline Ingestion pipeline

# Offline Ingestion pipeline



Partitioned Analyzed Daily Tweet Data

Earlybird

Earlybird

*earlybird machines load preprocessed data from HDFS for **indexing***

Earlybird segments on HDFS

# Offline Ingestion pipeline



like the realtime cluster they **share** complete segments via **HDFS**

Partitioned Analyzed Daily Tweet Data

Earlybird

Earlybird

...

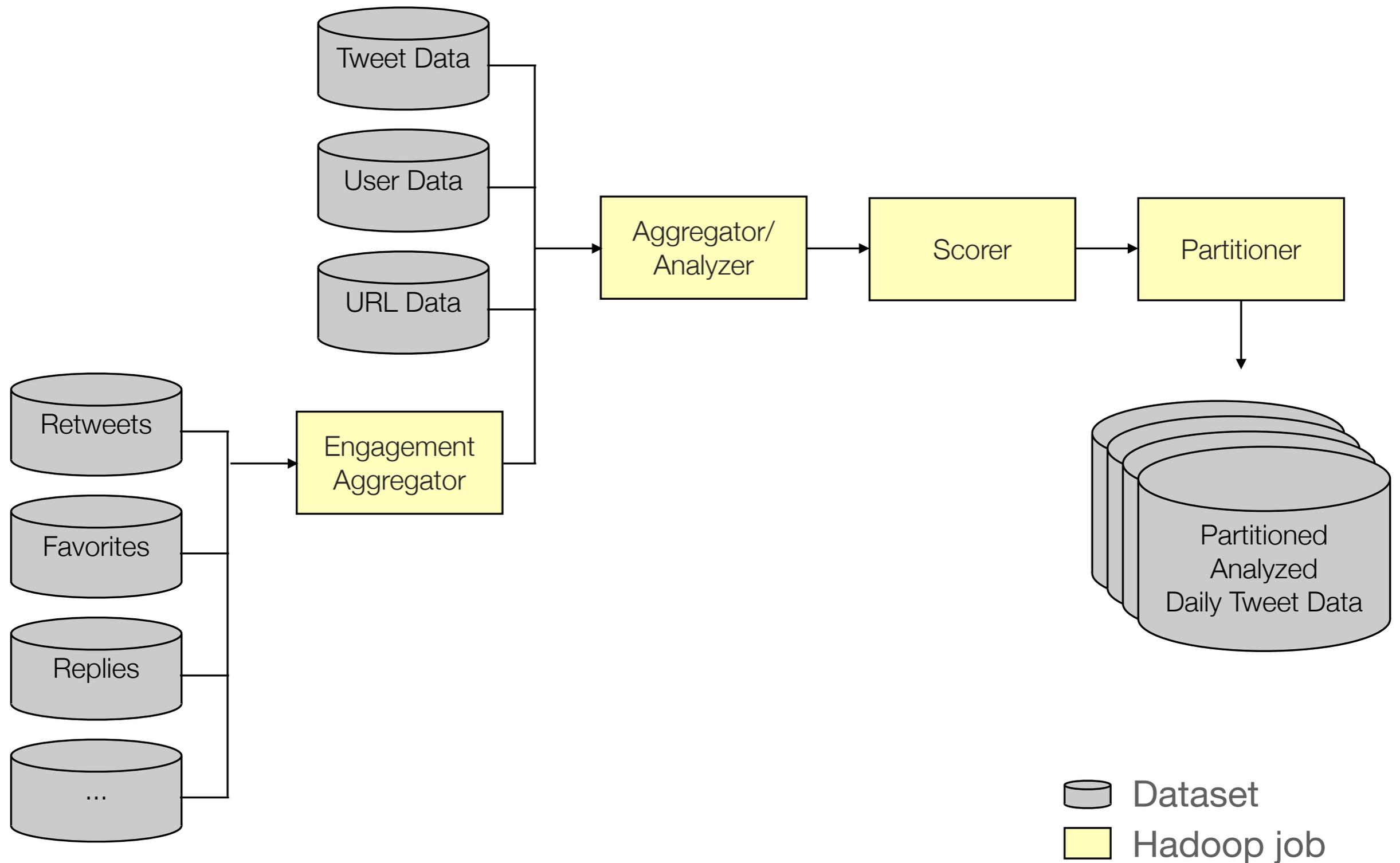Earlybird

Earlybird segments on HDFS

# 2013

## 2013

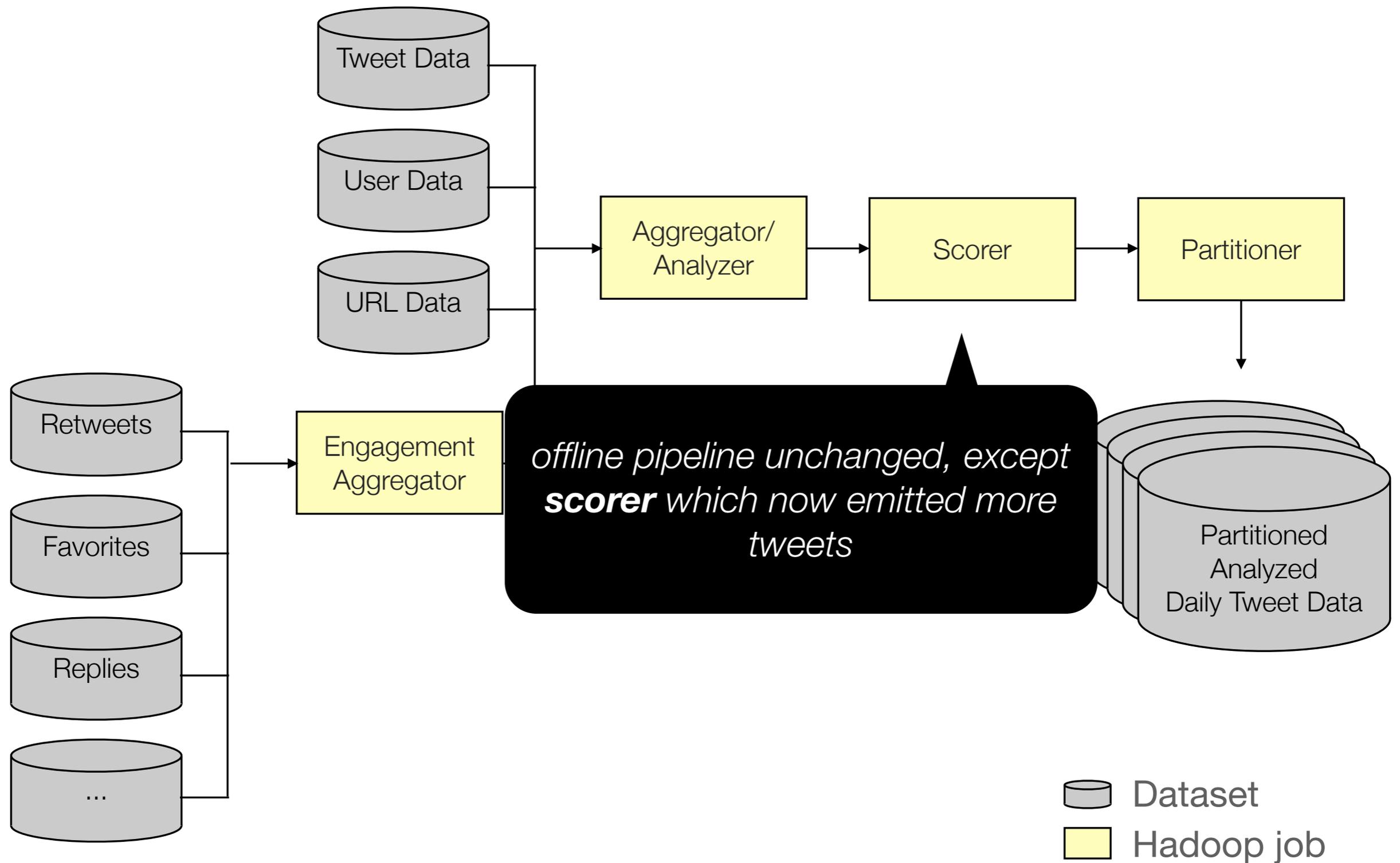Twitter launches **SSD-based historical index.**

# SSD-based historical index

- SSD-based index containing an order of magnitude more tweets than in-memory historical index

- Vanilla Lucene 4.x index format

- Explicit caching of forward indexes

- Rigorous hardware tuning for optimized IOPs

- "Traditional" hash-partitioned cluster layout
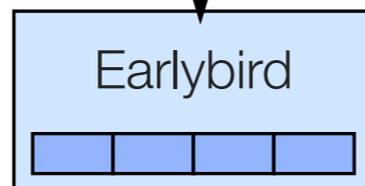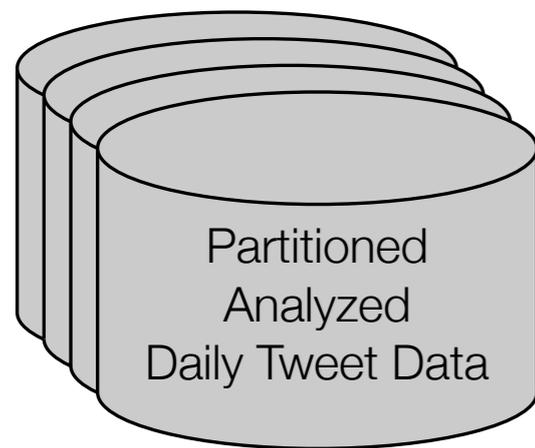
# Offline Ingestion pipeline
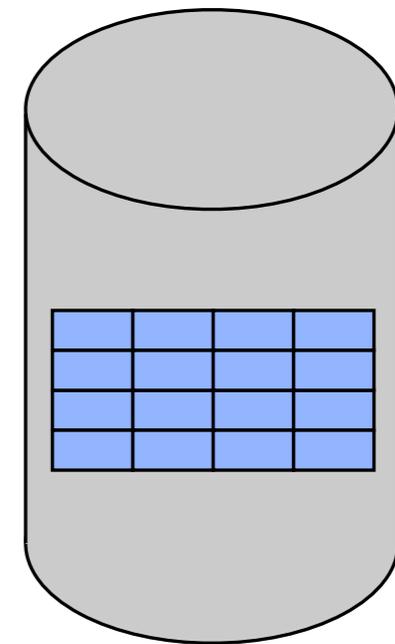
# Offline Ingestion pipeline

# Performance tuning

- Optimized hardware configuration for maximum SSD utilization

- Explicit packing and caching of Lucene DocValues for scoring

- Auxiliary skip lists using static (query independent) relevance signals

- Graceful degradation in disaster scenarios

2014

Twitter launches **Full Tweet Index.**

# Full Tweet Index

- SSD-based index containing all Tweets ever published

- An order of magnitude bigger index - new ways of scaling necessary:

    - Faster index builds

    - Expanding the index without repartitioning
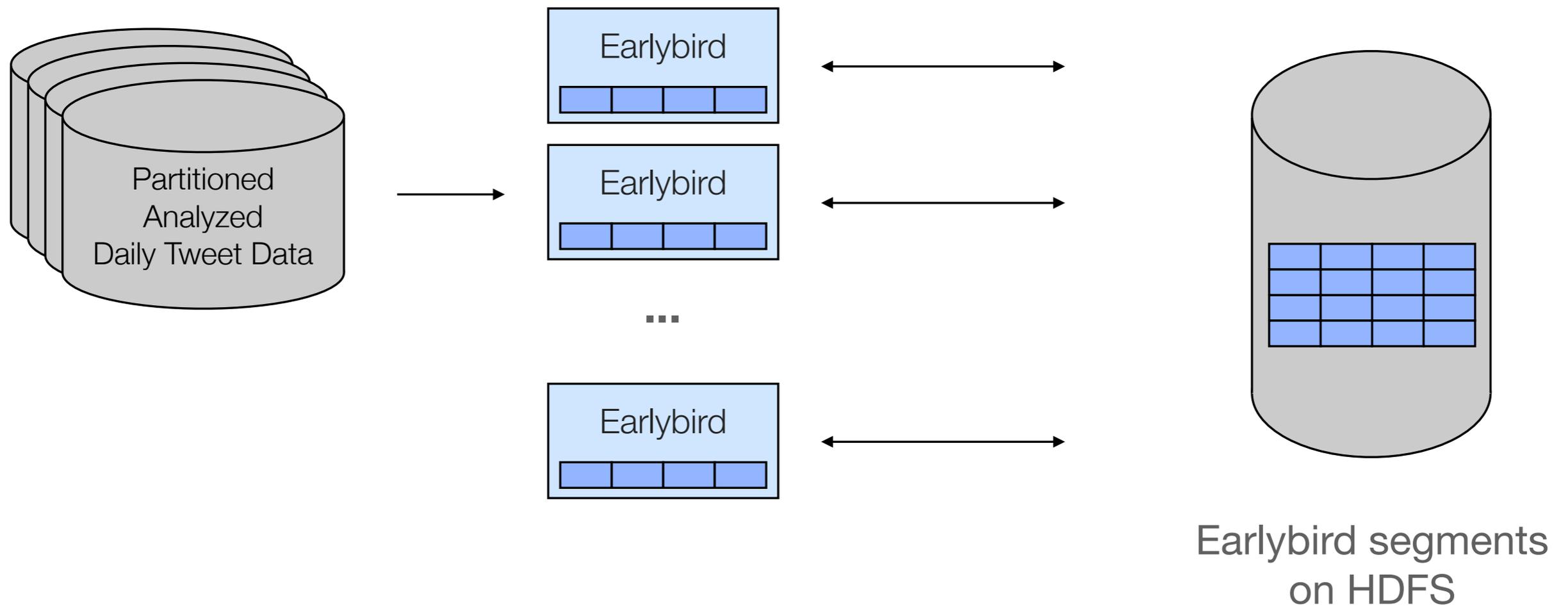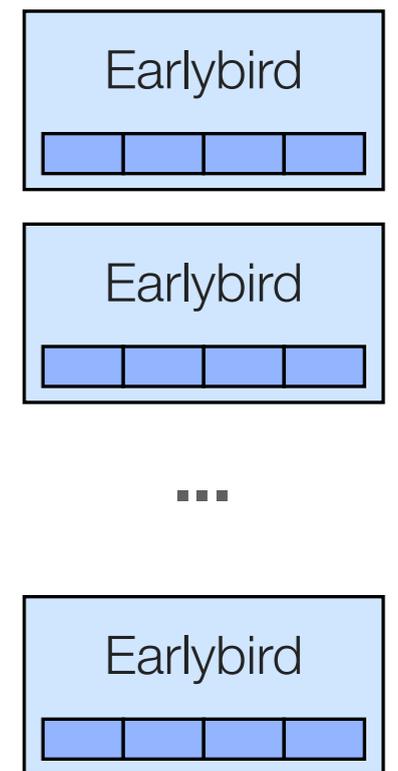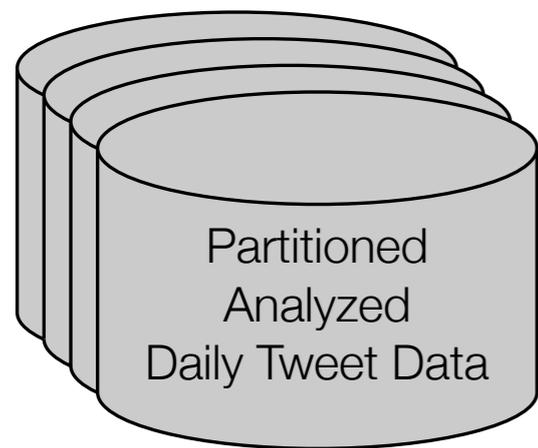
# Full Tweet Index

- SSD-based index containing all Tweets ever published

- An order of magnitude bigger index - new ways of scaling necessary:

  - **Faster index builds**
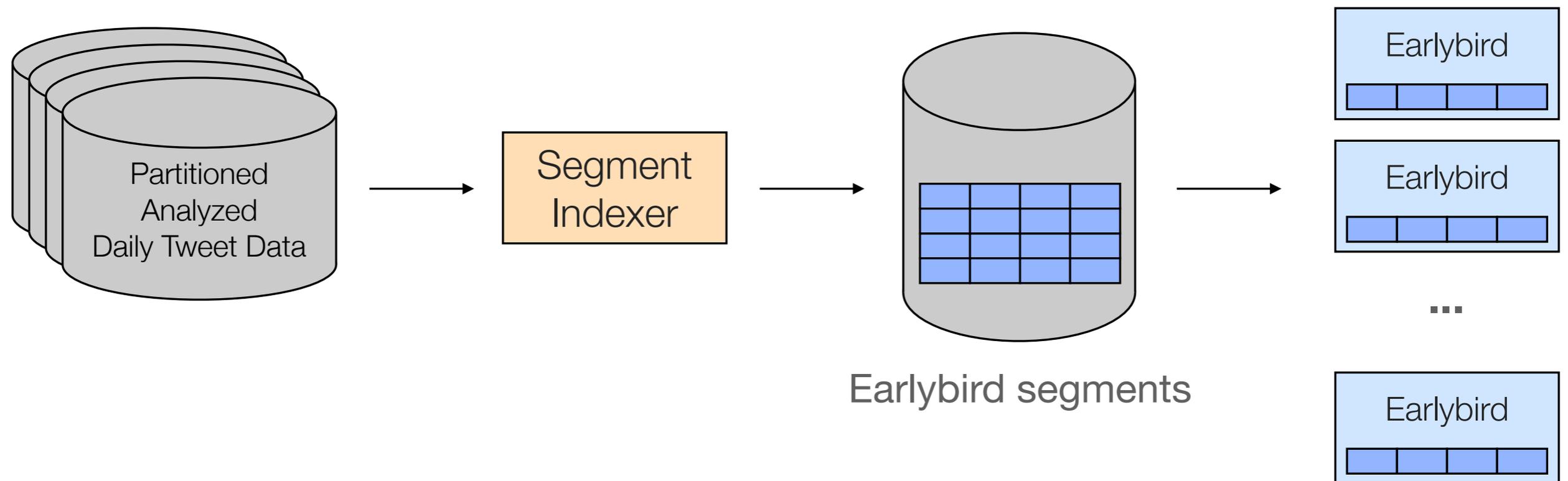
  - Expanding the index without repartitioning

# Offline Ingestion pipeline



Partitioned Analyzed Daily Tweet Data

Earlybird

Earlybird

...

Earlybird

Earlybird segments on HDFS

# Offline Ingestion pipeline

# Offline Ingestion pipeline



Partitioned Analyzed Daily Tweet Data → Segment Indexer → Earlybird segments → Earlybird ... Earlybird

Mesos job

# Offline Ingestion pipeline

massively parallelizable **mesos** jobs **invert** preprocessed data and write **earlybird segments**

Partitioned Analyzed Daily Tweet Data

Segment Indexer

Earlybird segments

Earlybird

Earlybird

...

Earlybird

Mesos job

# Offline Ingestion pipeline



Partitioned Analyzed Daily Tweet Data → Segment Indexer → Earlybird segments

*a segment consists of **multiple days** of tweets; max segment size enforced*

Earlybird
Earlybird
...
Earlybird

Mesos job

# Offline Ingestion pipeline



Partitioned Analyzed Daily Tweet Data → Segment Indexer → [HDFS] → Earlybird ... Earlybird

*earlybird machines **copy** segments from HDFS for **serving***

Mesos job

# Offline Ingestion pipeline



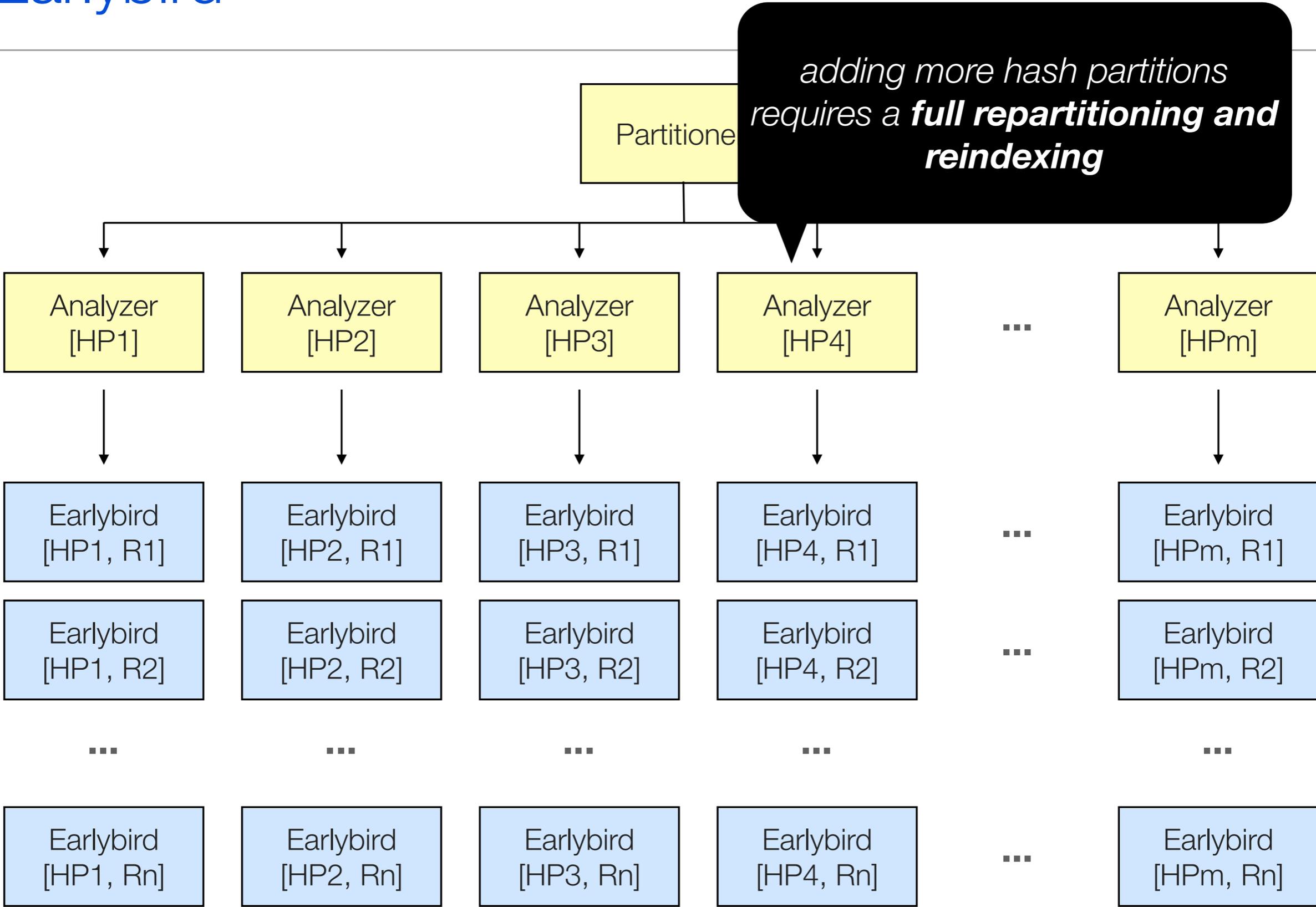Partitioned Analyzed Daily Tweet Data → Segment Indexer → Earlybird segments → Earlybird ... Earlybird

segment indexer produces **vanilla lucene 4.x** segments
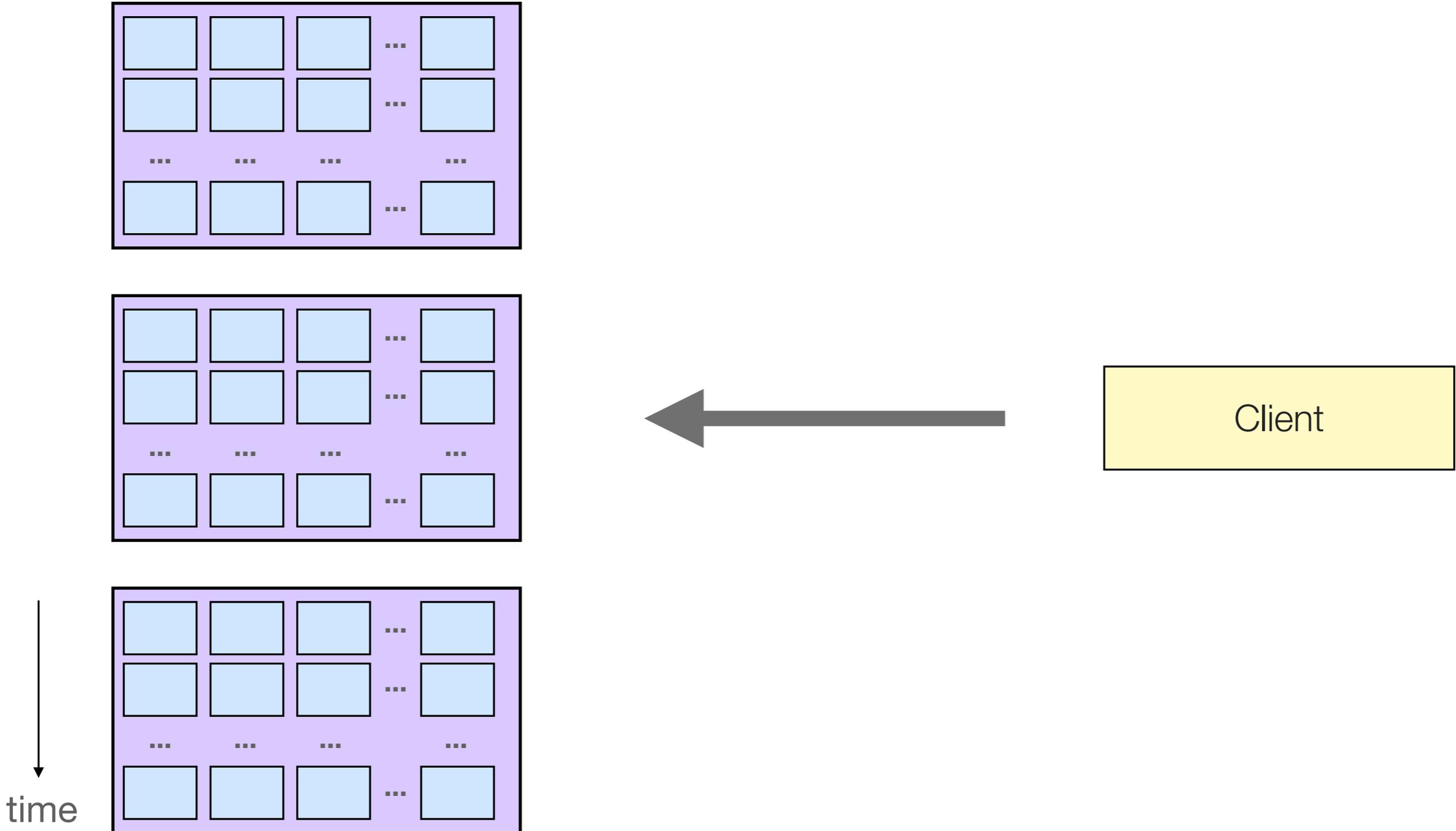
Mesos job

# Full Tweet Index

- SSD-based index containing all Tweets ever published

- An order of magnitude bigger index - new ways of scaling necessary:

  - Faster index builds

  - **Expanding the index without repartitioning**
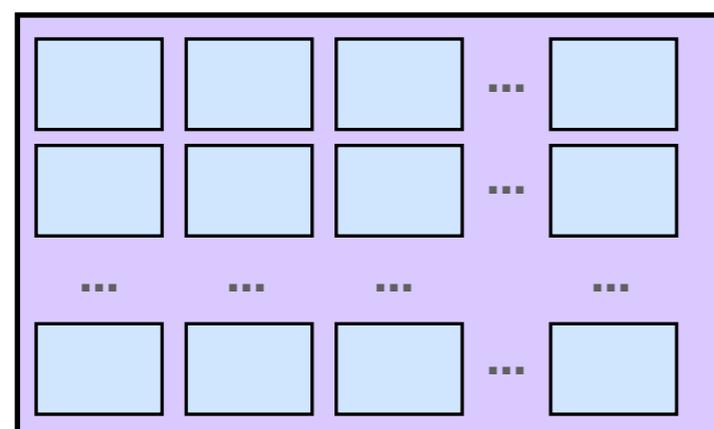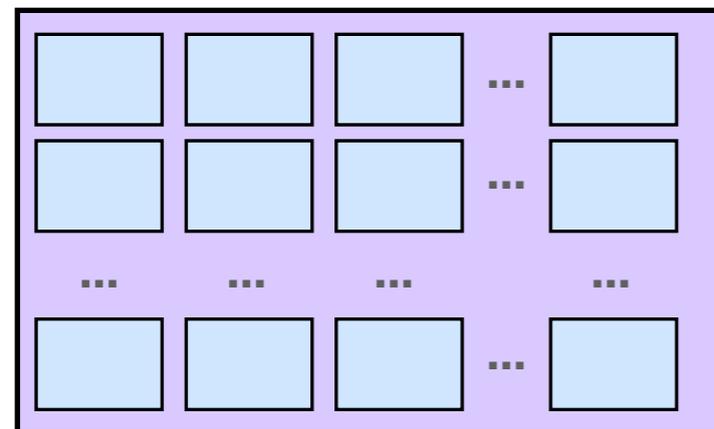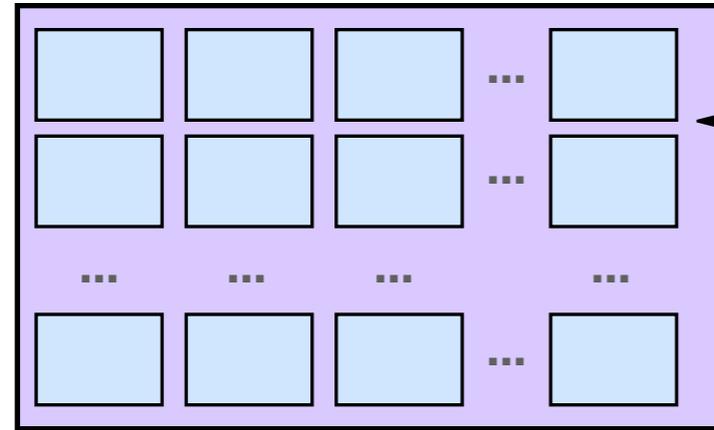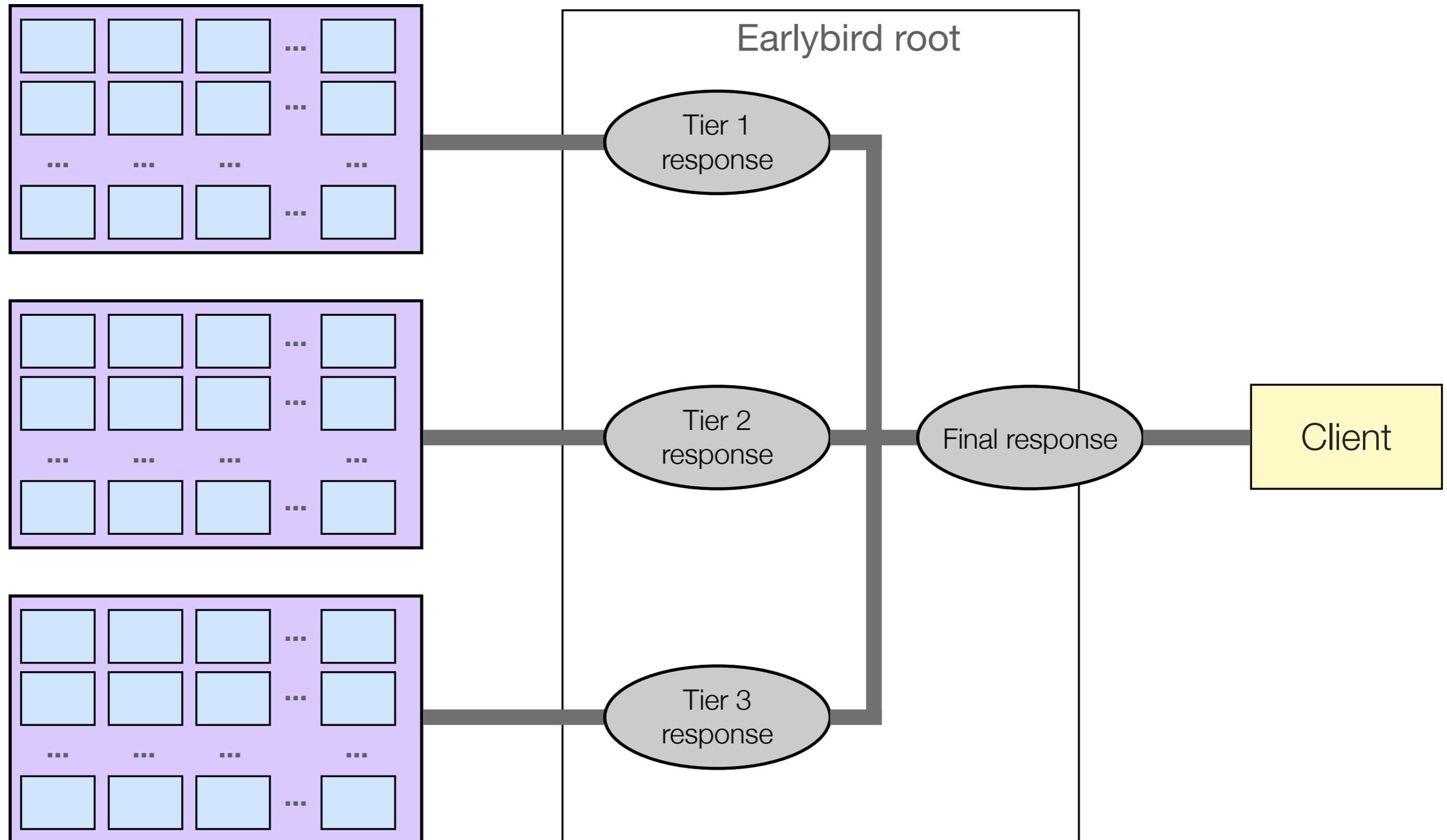
# Earlybird

# Time Tiers



time

# Time Tiers



each **tier** covers a **fixed time range**; each tier has a fixed number of hash partitions and replicas

Client

time
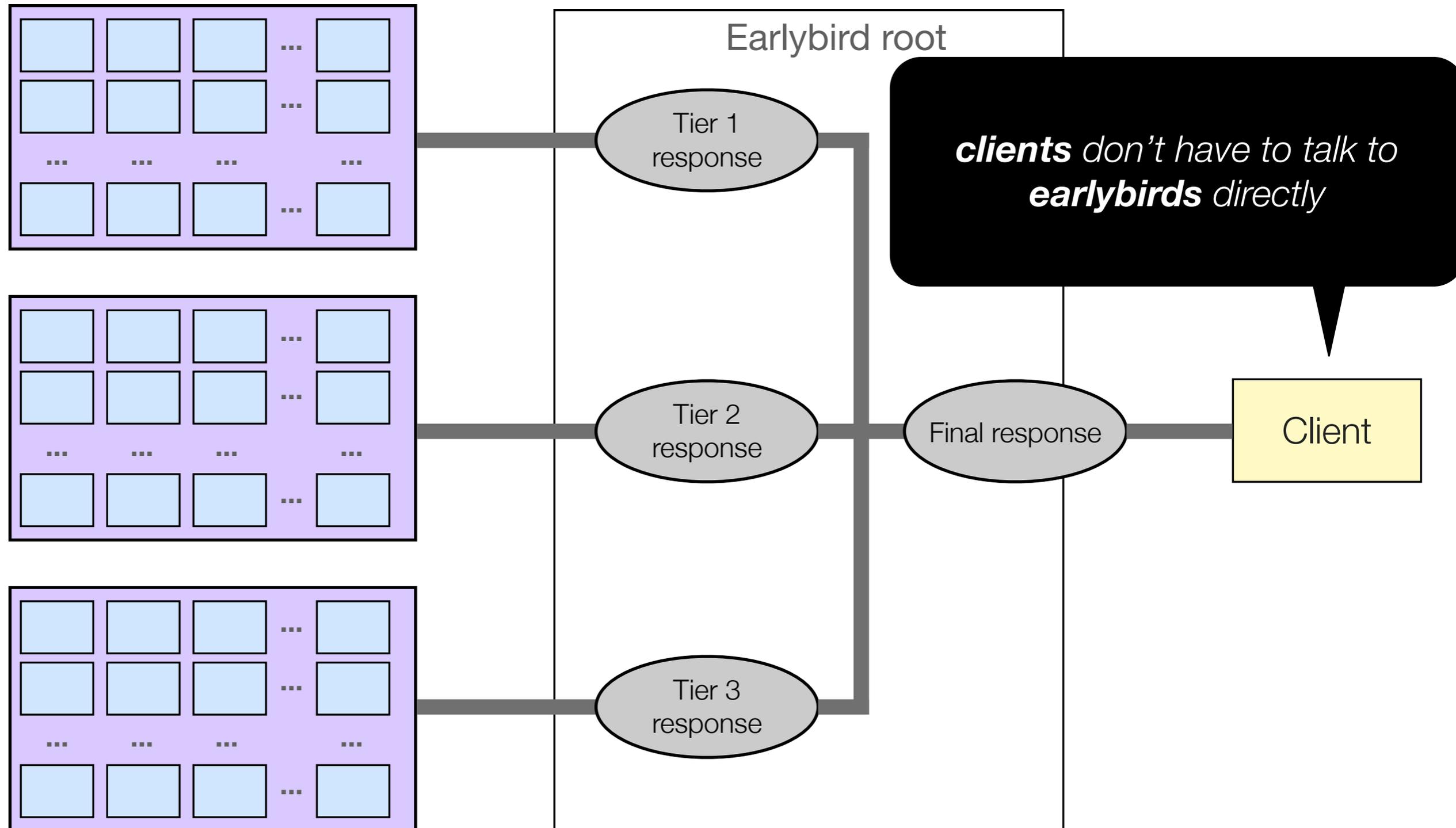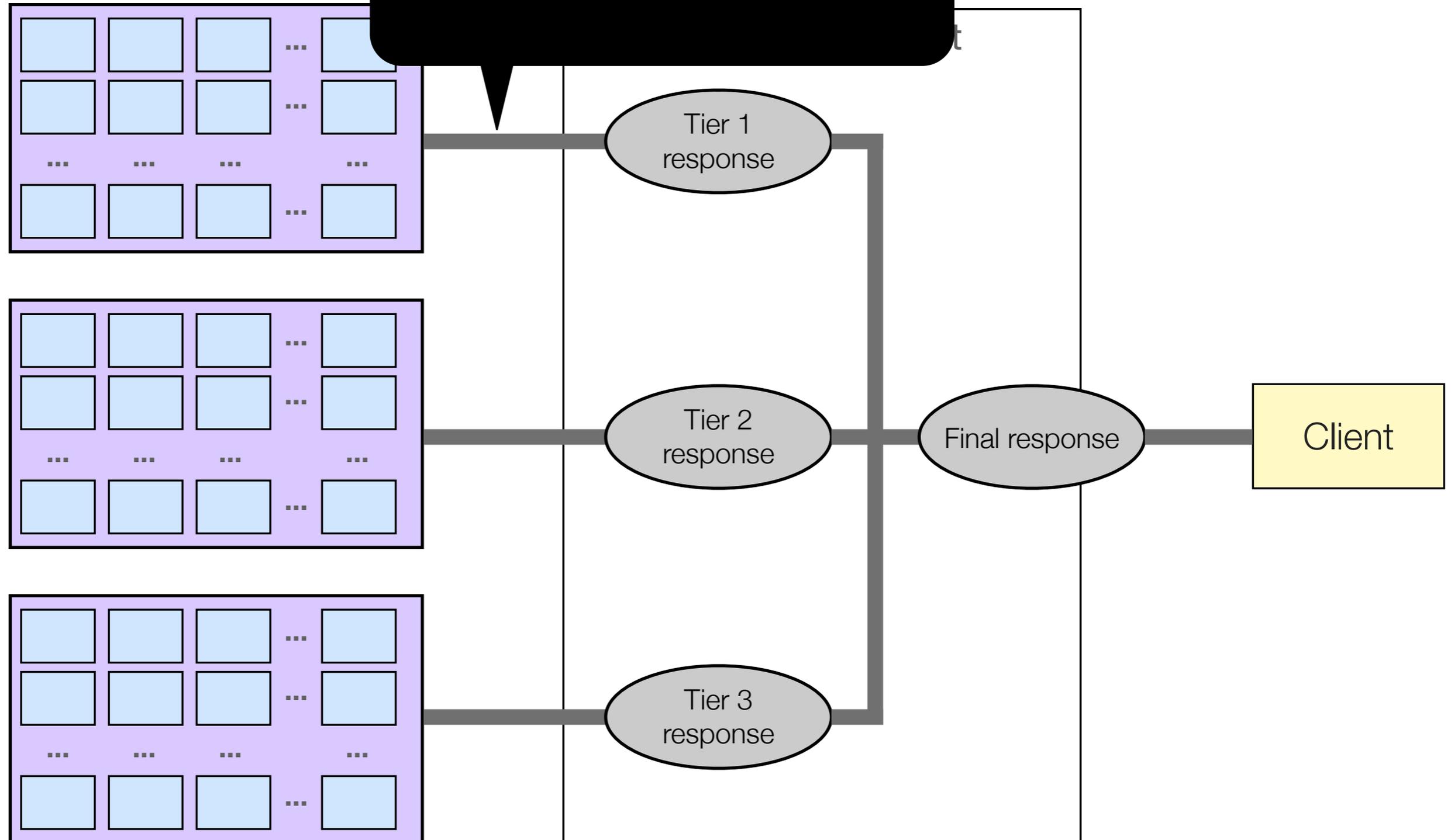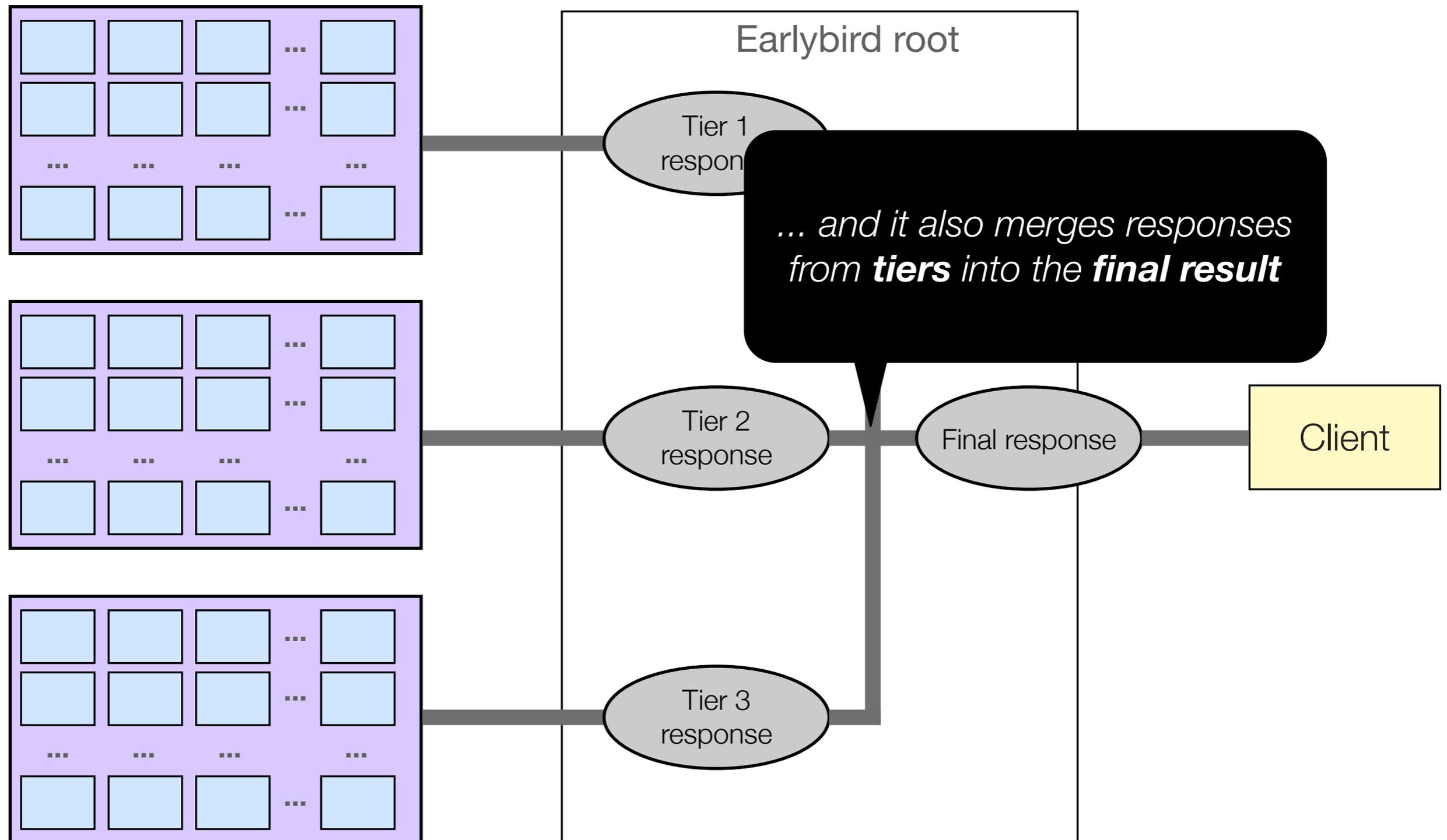
# Search roots

# Search roots

# Search roots



earlybird root merges results from multiple hash partitions …

time

Tier 1 response

Tier 2 response

Tier 3 response

Final response

Client

# Search roots

# Full Tweet Index

- SSD-based index containing all Tweets ever published

- Data structures to make previous historical indexes obsolete

- New tiered cluster layout for easy index scaling

- New Mesos-based index builder system

# In-memory Real-time Index

- Highly optimized for GC - all data is stored in blocked native arrays

- v1: Optimized for tweets with a term position limit of 255

- v2: Support for 32 bit positions without performance degradation

- v2: Basic support for out-of-order posting list inserts



**Michael Busch** 🐦 Follow
@michibusch

Removing code written 4 years ago. You helped people a trillion times find what they were searching for. #RIP #makingroomfornewthings

1:22 PM - 21 Jan 2014 from Northwest Marin, CA, United States

4 FAVORITES

# 2015

Outlook

# Outlook

- **Parallel indexing pipelines for faster index manipulation**

- Domain-independent core indexing library that combines real-time and offline index technology

# Modifying indexes

- The tiered architecture allows us to grow the index without having to reindex existing data

- But appending new fields to existing indexes or e.g. upgrading analyzers requires a full reprocessing of all data

- Idea: Introduce segment slices and parallel ingestion pipelines
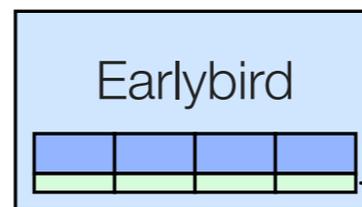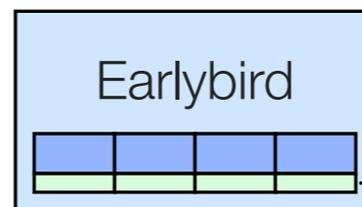
# Segment Slices

Earlybird

# Segment Slices



Earlybird

we want to append a **new field** to an existing segment; however lucene segments are **immutable**
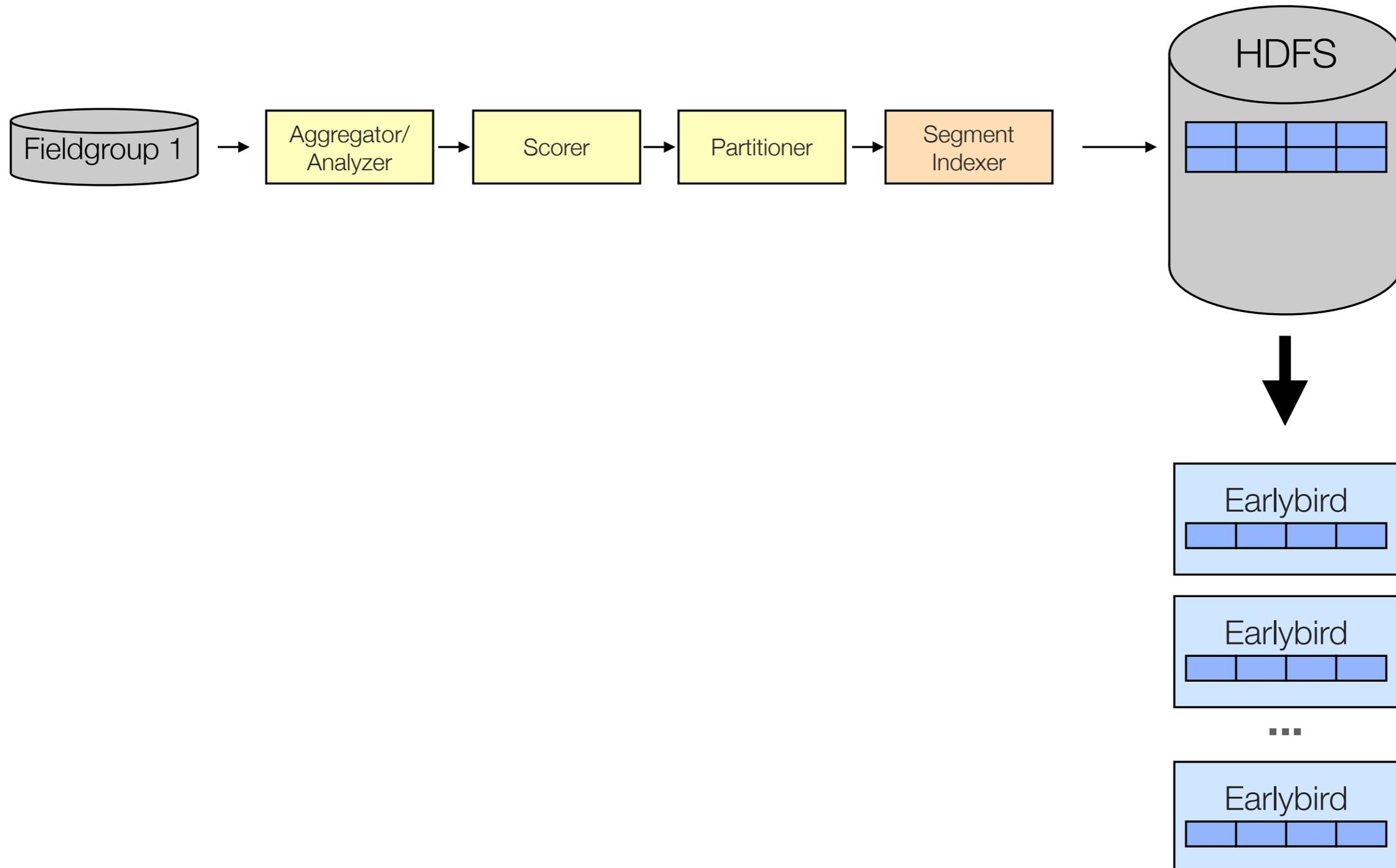
# Segment Slices

Earlybird

*index a separate parallel segment covering the **same doc id range***
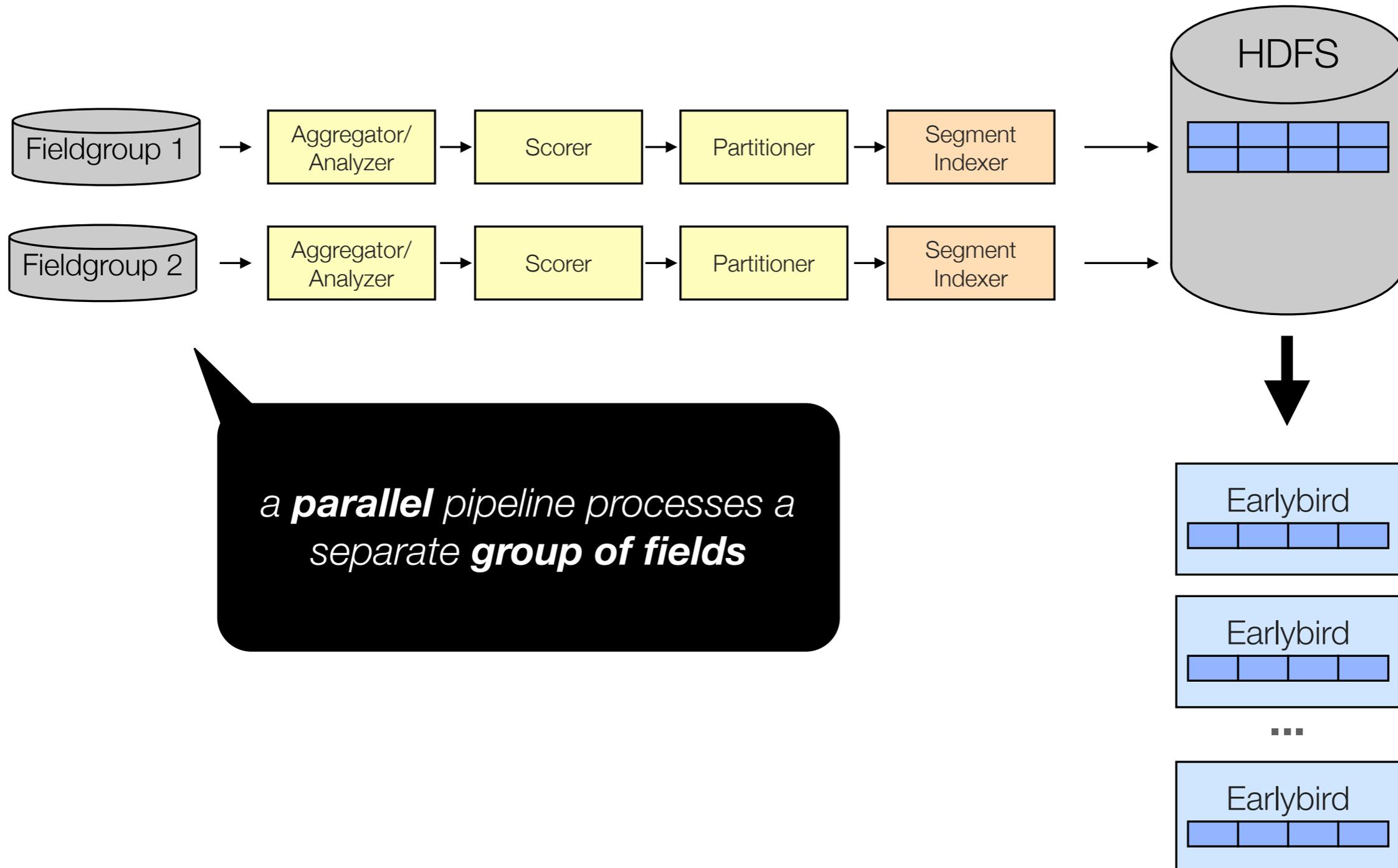
# Segment Slices

Earlybird

*using lucene's **ParallelIndexReader** we can read both segment slides as if they were a single one*
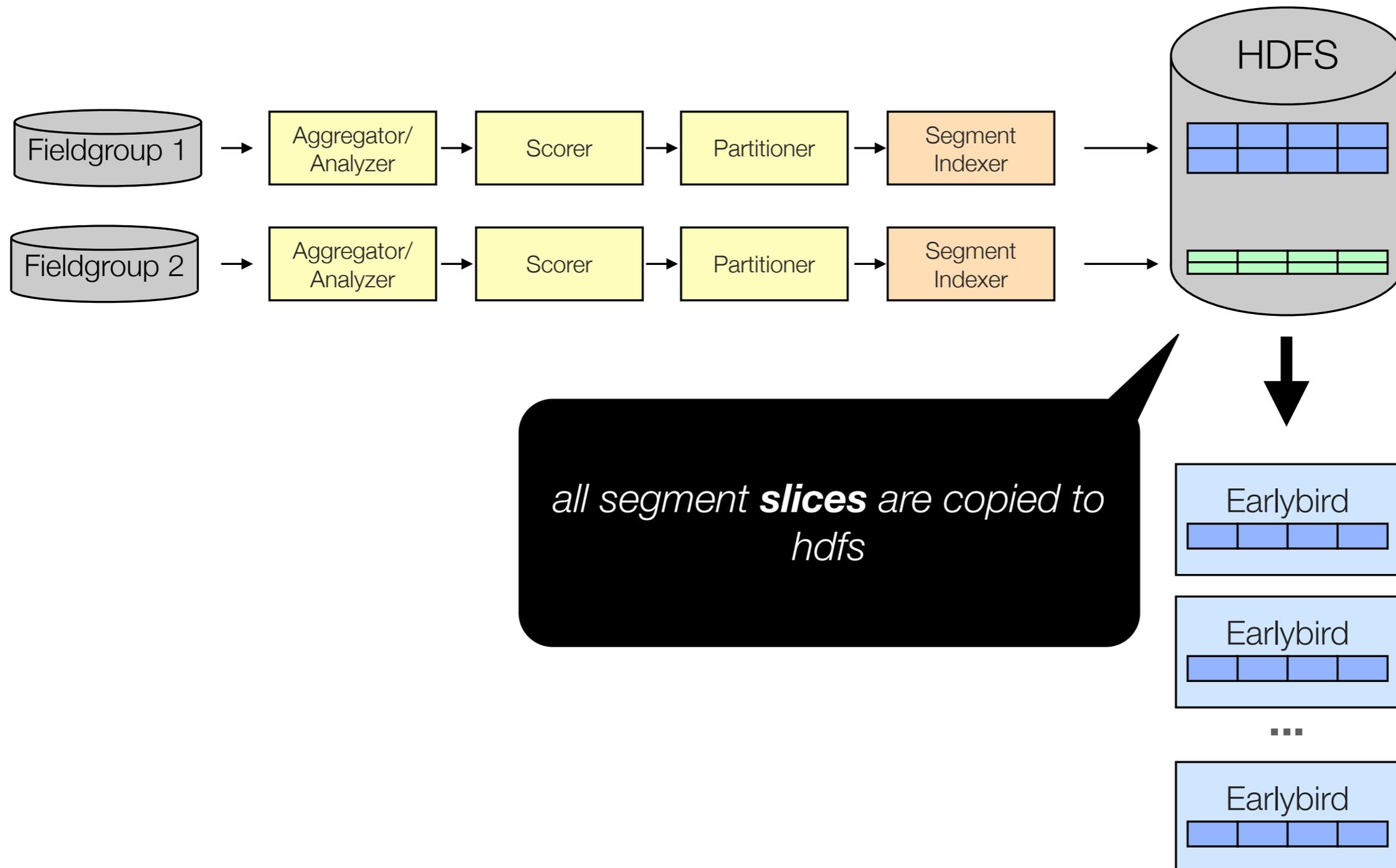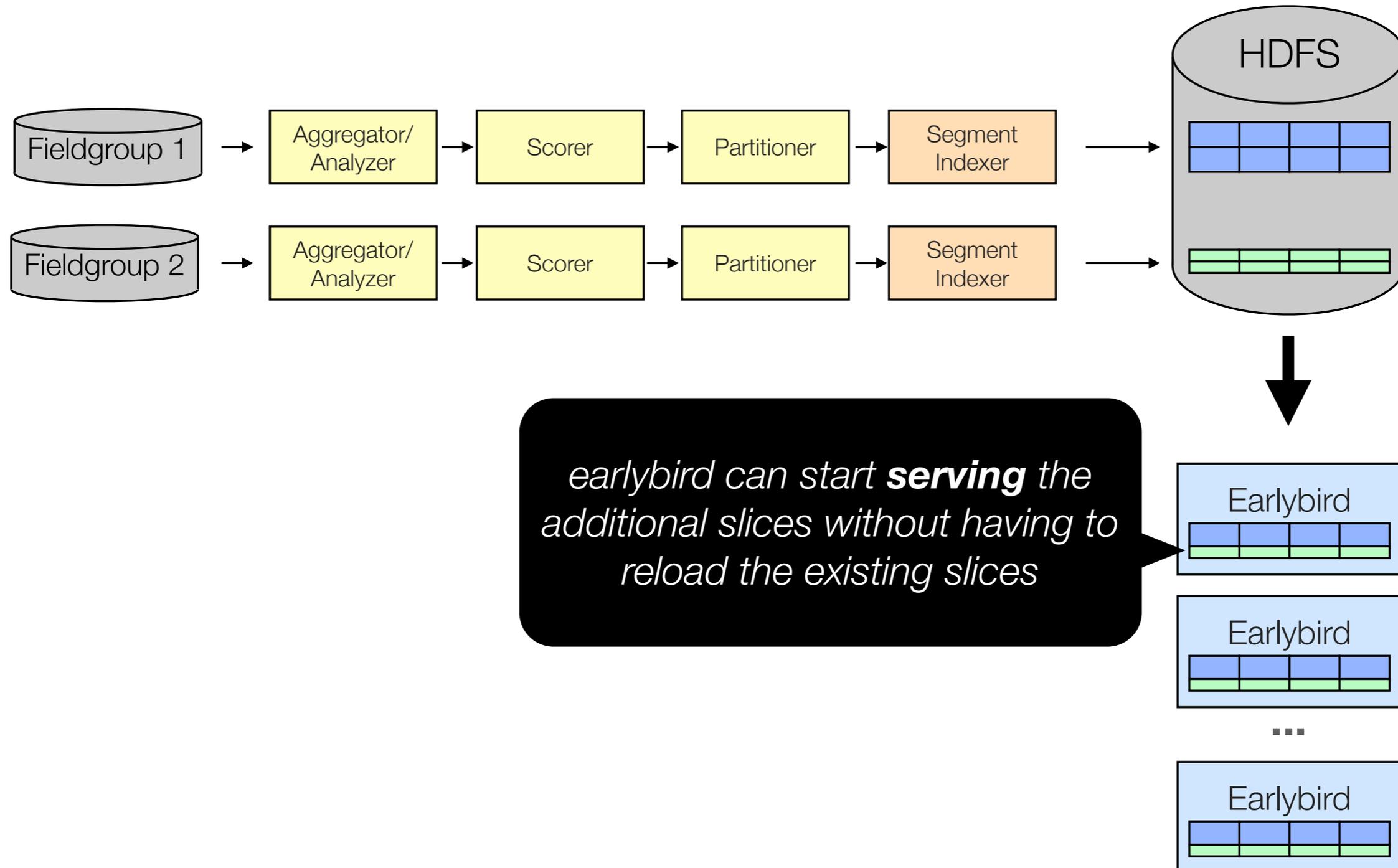
# Parallel Ingestion pipeline

# Parallel Ingestion pipeline
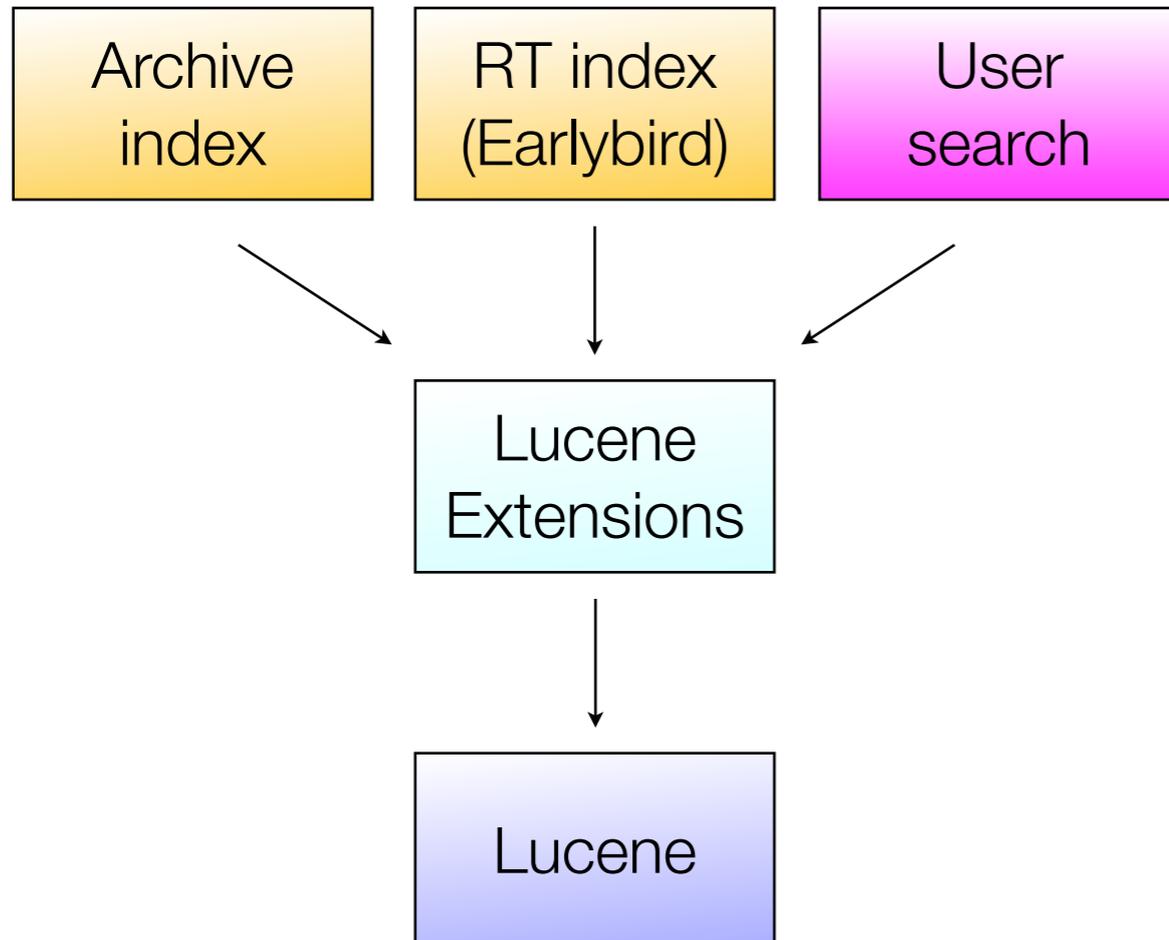
# Parallel Ingestion pipeline

# Parallel Ingestion pipeline

# Outlook

- Parallel indexing pipelines for faster index manipulation

- **Domain-independent core indexing library that combines real-time and offline index technology**

# Search Architecture



- New Lucene extension package

- This package is truly generic and has no dependency on an actual product/index

- It contains Twitter's extensions for real-time search, a thin segment management layer and other features

# Lucene Extension Library

- Abstraction layer for Lucene index segments

- Real-time writer for in-memory index segments

- Schema-based Lucene document factory

- Real-time faceting

# Lucene Extension Library

- API layer for Lucene segments

    - *IndexSegmentWriter

    - *IndexSegmentAtomicReader

- Two implementations

    - In-memory: RealtimeIndexSegmentWriter (and reader)

    - On-disk: LuceneIndexSegmentWriter (and reader)

# Lucene Extension Library

- IndexSegments can be built ...

  - in realtime

  - on Mesos or Hadoop (Mapreduce)

  - locally on serving machines

- Cluster-management code that deals with IndexSegments

  - Share segments across serving machines using HDFS

  - Can rebuild segments (e.g. to upgrade Lucene version, change data schema, etc.)

# Demo