# Stream Processing beyond streaming data – Batch, Streaming, and Applications

Timo Walther

PMC of Apache Flink

Ververica (formerly dataArtisans, now part of Alibaba Group)

Follow me: @twalthr (yes, without the e)

# Alternative Talk Titles

"Batch is a special case of streaming"

"If all you have is a Squirrel, everything looks like a stream"

"What's taking you so long to
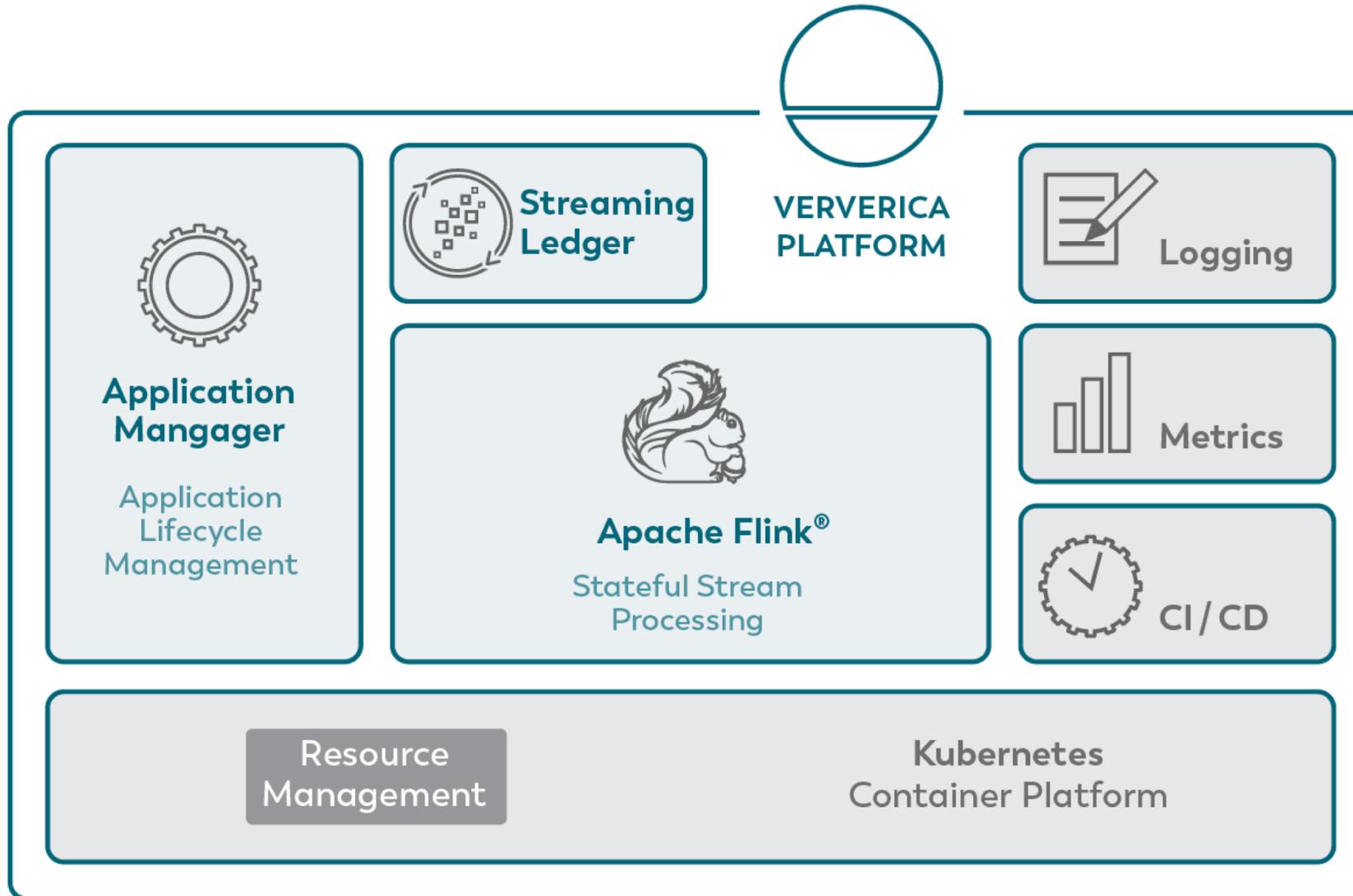merge DataSet and DataStream?"

# About Ververica

Original creators of
Apache Flink®

Enterprise
Stream Processing

# Ververica Platform

# Flink 101

# Some Apache Flink Users

# Stream Processing with Apache Flink at Alibaba

## The "Singles Day" (11/11)

| machines | queries | throughput | latency | state size |
|----------|---------|------------|---------|------------|
| 10K | 10K | 1.7B events / sec | Sub-Second | 100TB |

# Apache Flink

## Stateful Computations over Data Streams

# The Flink Vision

| Batch Processing | Continuous Processing | Data Pipelines | Streaming Analytics | Event-driven Applications | Transactional Applications |
|---|---|---|---|---|---|

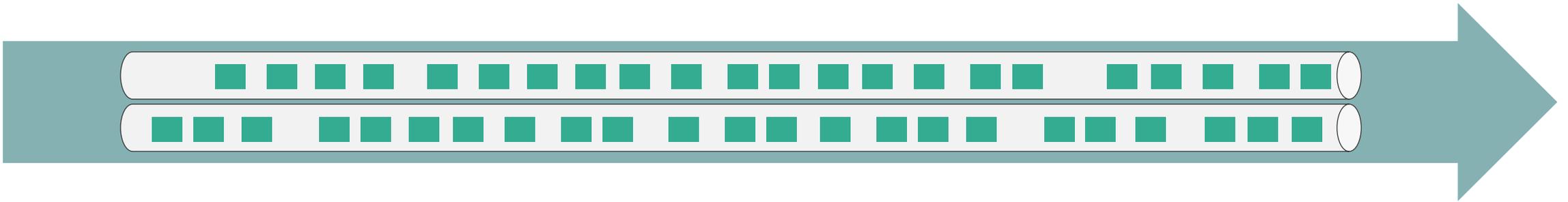more lag time                                                                    more real time

Unifying data processing based on thinking in data streams

And building the first open source system to cover that spectrum
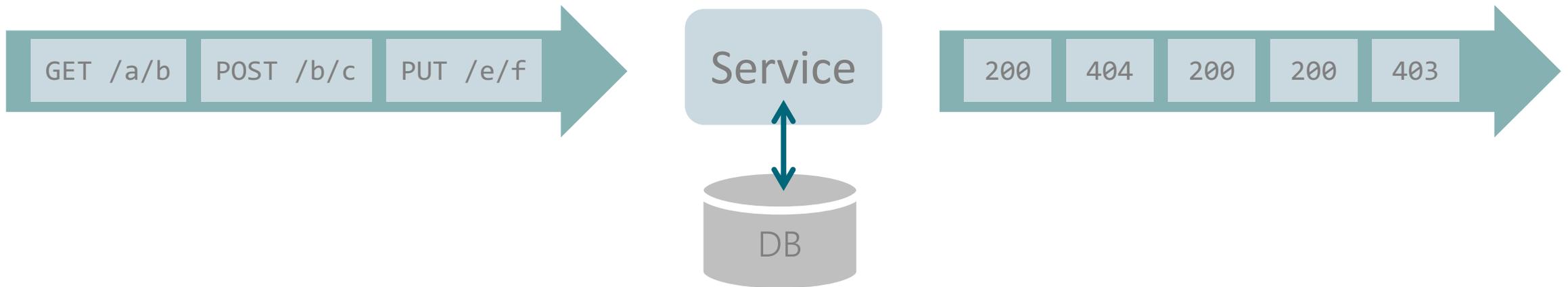
# Everything is a Stream

Streams Of Records in a Log or MQ
[e.g., Apache Kafka or AWS Kinesis ...]

# Everything is a Stream

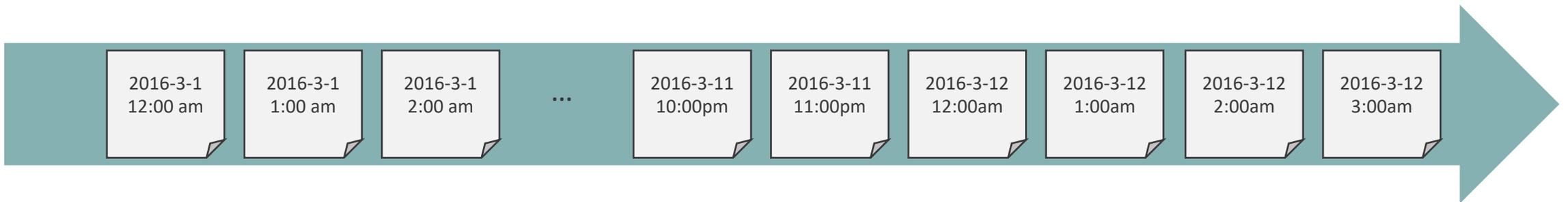Stream of Requests/Responses to/from Services

| GET /a/b | POST /b/c | PUT /e/f |
|----------|-----------|----------|

Service

DB

| 200 | 404 | 200 | 200 | 403 |
|-----|-----|-----|-----|-----|

→ event sourcing architecture

# Everything is a Stream

## Stream of Rows in a Table or in Files



| 2016-3-1 12:00 am | 2016-3-1 1:00 am | 2016-3-1 2:00 am | ... | 2016-3-11 10:00pm | 2016-3-11 11:00pm | 2016-3-12 12:00am | 2016-3-12 1:00am | 2016-3-12 2:00am | 2016-3-12 3:00am |

# Everything is a Stream

Stream of Rows in a Table or in Files

| 2016-3-1 12:00 am | 2016-3-1 1:00 am | 2016-3-1 2:00 am | ... | 2016-3-11 10:00pm | 2016-3-11 11:00pm | 2016-3-12 12:00am | 2016-3-12 1:00am | 2016-3-12 2:00am | 2016-3-12 3:00am |

a batch

# Everything is a Stream

Streams may span storage systems

Parquet files

Avro records

| 2016-3-1 12:00 am | 2016-3-1 1:00 am | 2016-3-1 2:00 am | ... | 2016-3-11 10:00pm | 2016-3-11 11:00pm |

more distant past

*(e.g., compressed files in DFS/Object Store)*
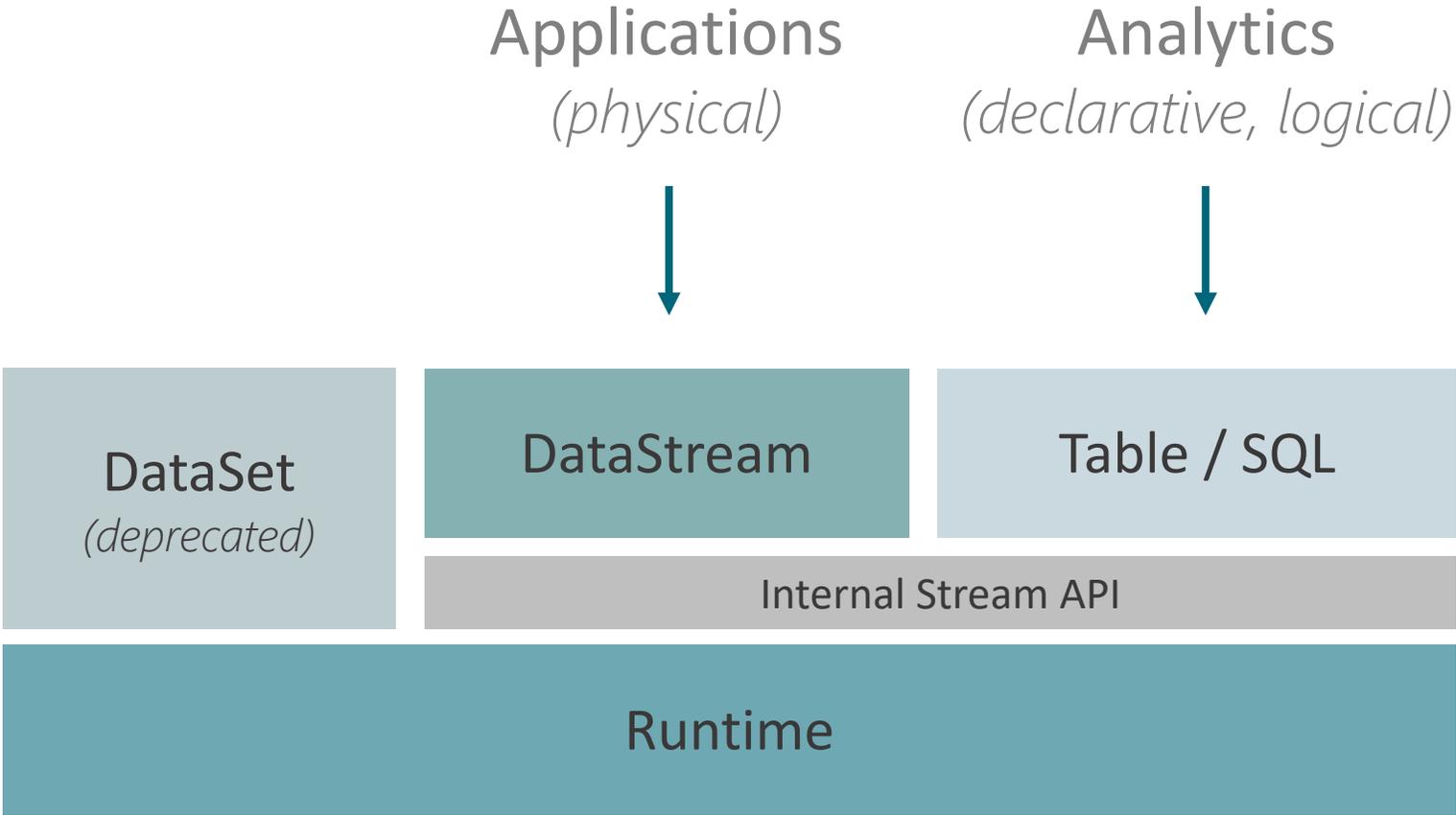
recent past

*(e.g., events in MQ/Log)*

Easy way to bootstrap an application with past data and then let it continue with real time data

# APIs

# APIs to Support these Use Cases

**Applications**
*(physical)*

**Analytics**
*(declarative, logical)*

| DataSet *(deprecated)* | DataStream | Table / SQL |
|---|---|---|

Internal Stream API

Runtime

# DataStream API (Functional Java / Scala)

```scala
val lines: DataStream[String] = env.addSource(new FlinkKafkaConsumer011(…))
```
Source

```scala
val events: DataStream[Event] = lines.map((line) => parse(line))
```
Transformation

```scala
val stats: DataStream[Statistic] = stream
    .keyBy("sensor")
    .timeWindow(Time.seconds(5))
    .sum(new MyAggregationFunction())
```
Windowed Transformation

```scala
stats.addSink(new StreamingFileSink(path))
```
Sink

# DataStream API Process Functions

```java
public void processElement1(Transaction txn, Context ctx, Collector<Transaction> out) {
    // keep the transaction in the internal state until the approval comes
    pendingTransaction.update(txn);
    // schedule a timer to trigger the timeout
    ctx.timerService().registerProcessingTimeTimer(txn.getTimestamp() + TIMEOUT_MILLIS);
}

public void processElement2(ApproveOrReject approval, Context ctx, Collector<Transaction> out) {
    // get and remove the transaction from the state
    Transaction txn = pendingTransaction.value();
    pendingTransaction.clear();
    // forward the transaction to the main stream
    out.collect(txn);
}

public void onTimer(long timestamp, OnTimerContext ctx, Collector<Transaction> out) {
    // check if the transaction is still there, in which case it would be timed out
    Transaction txn = pendingTransaction.value();
    if (txn != null) {
        // write to the timeout stream
        ctx.output(TIMEOUT_STREAM, txn);
        pendingTransaction.clear();
    }
}
```

# SQL / Table API – Batch style (fix data set as input)
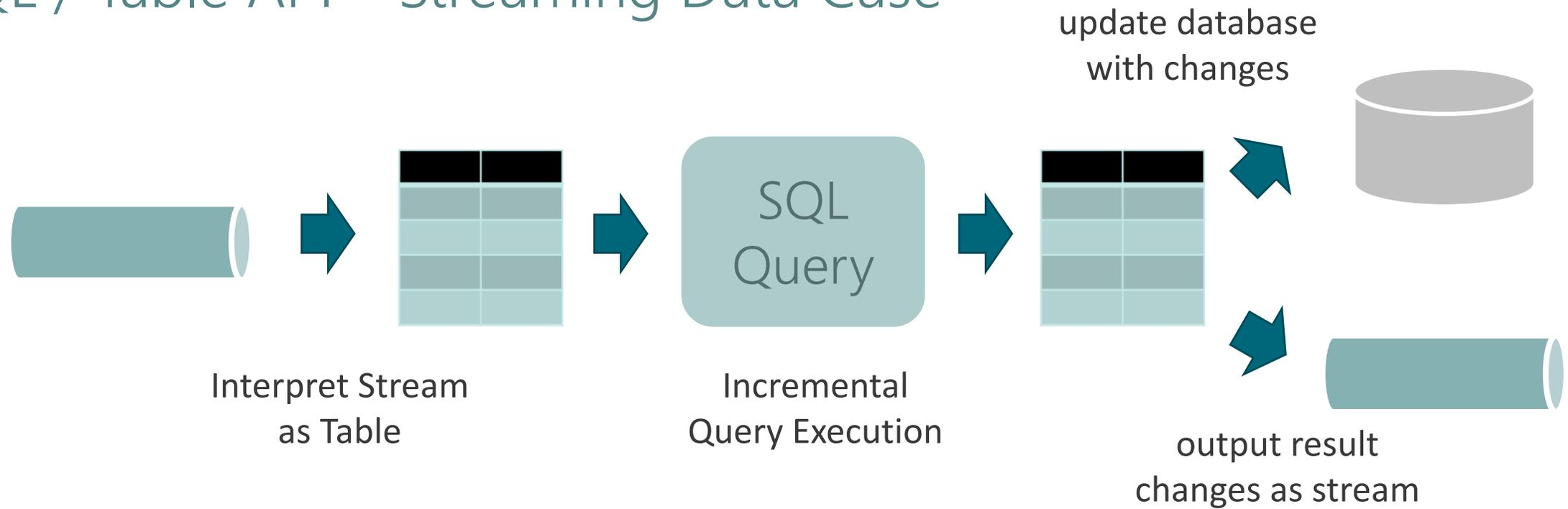


Batch Query
Execution

Full TPC-H support in
Flink 1.9 with Blink query engine

Full TPC-DS support
targeted for Flink 1.10

```
SELECT
    room,
    TUMBLE_END(rowtime, INTERVAL '1' HOUR),
    AVG(temperature)
FROM
    sensors
GROUP BY
    TUMBLE(rowtime, INTERVAL '1' HOUR), room
```
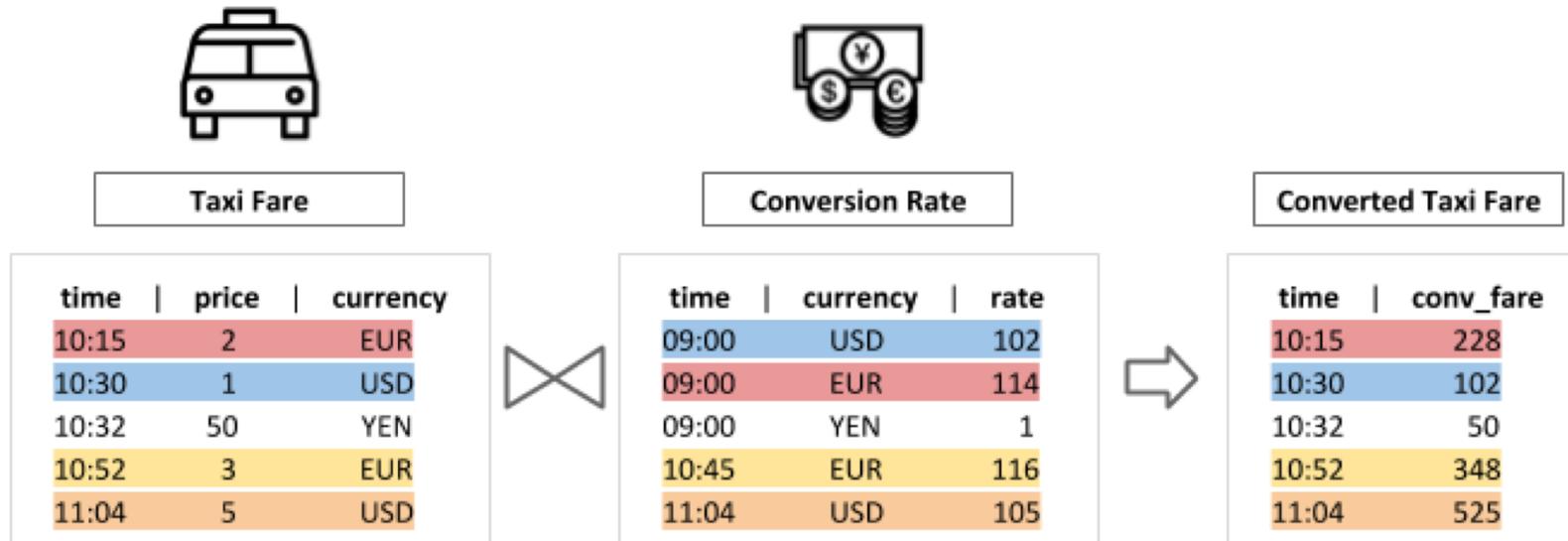
# SQL / Table API – Streaming Data Case



Interpret Stream as Table

Incremental Query Execution

update database with changes

output result changes as stream

```sql
SELECT
    room,
    TUMBLE_END(rowtime, INTERVAL '1' HOUR),
    AVG(temperature)
FROM
    sensors
GROUP BY
    TUMBLE(rowtime, INTERVAL '1' HOUR), room
```

# SQL / Table API – Temporal Joins Example



| Taxi Fare | | |
| --- | --- | --- |
| time | price | currency |
| 10:15 | 2 | EUR |
| 10:30 | 1 | USD |
| 10:32 | 50 | YEN |
| 10:52 | 3 | EUR |
| 11:04 | 5 | USD |

| Conversion Rate | | |
| --- | --- | --- |
| time | currency | rate |
| 09:00 | USD | 102 |
| 09:00 | EUR | 114 |
| 09:00 | YEN | 1 |
| 10:45 | EUR | 116 |
| 11:04 | USD | 105 |

| Converted Taxi Fare | |
| --- | --- |
| time | conv_fare |
| 10:15 | 228 |
| 10:30 | 102 |
| 10:32 | 50 |
| 10:52 | 348 |
| 11:04 | 525 |

```
SELECT tf.time
       tf.price * rh.rate as conv_fare

FROM taxiFare AS tf

LATERAL TABLE (Rates(tf.time)) AS rh

WHERE tf.currency = rh.currency;
```

# SQL / Table API – Event Pattern Matching Example

```sql
SELECT rideId, timeDiff(startT, endT) / 60000 AS durationMin
FROM Rides
MATCH_RECOGNIZE (
  PARTITION BY rideId
  ORDER BY rideTime
  MEASURES
    S.rideTime AS startT,
    E.rideTime AS endT
  AFTER MATCH SKIP PAST LAST ROW
  PATTERN (S E)
  DEFINE
    S AS S.isStart,
    E AS NOT E.isStart
);
```

# The Relationship between Batch and Streaming
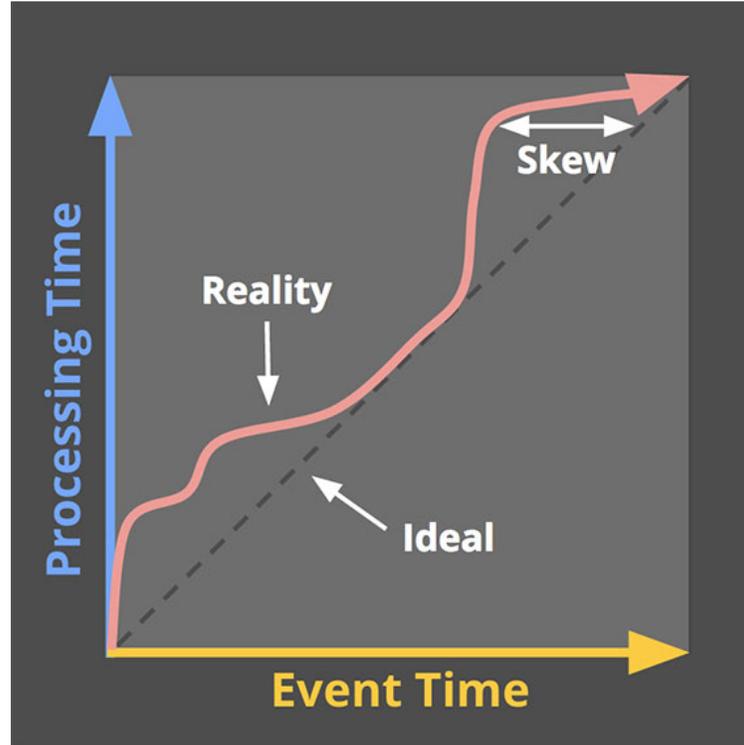
# Batch Processing is a special case of Stream Processing

A batch is just a bounded stream.



That is about 60% of the truth...

# The remaining 40% of the truth



... never seen this in Batch Processing, though.

## The (Event-time) Watermark

# The remaining 40% of the truth

## Continuous Streaming

Data is incomplete

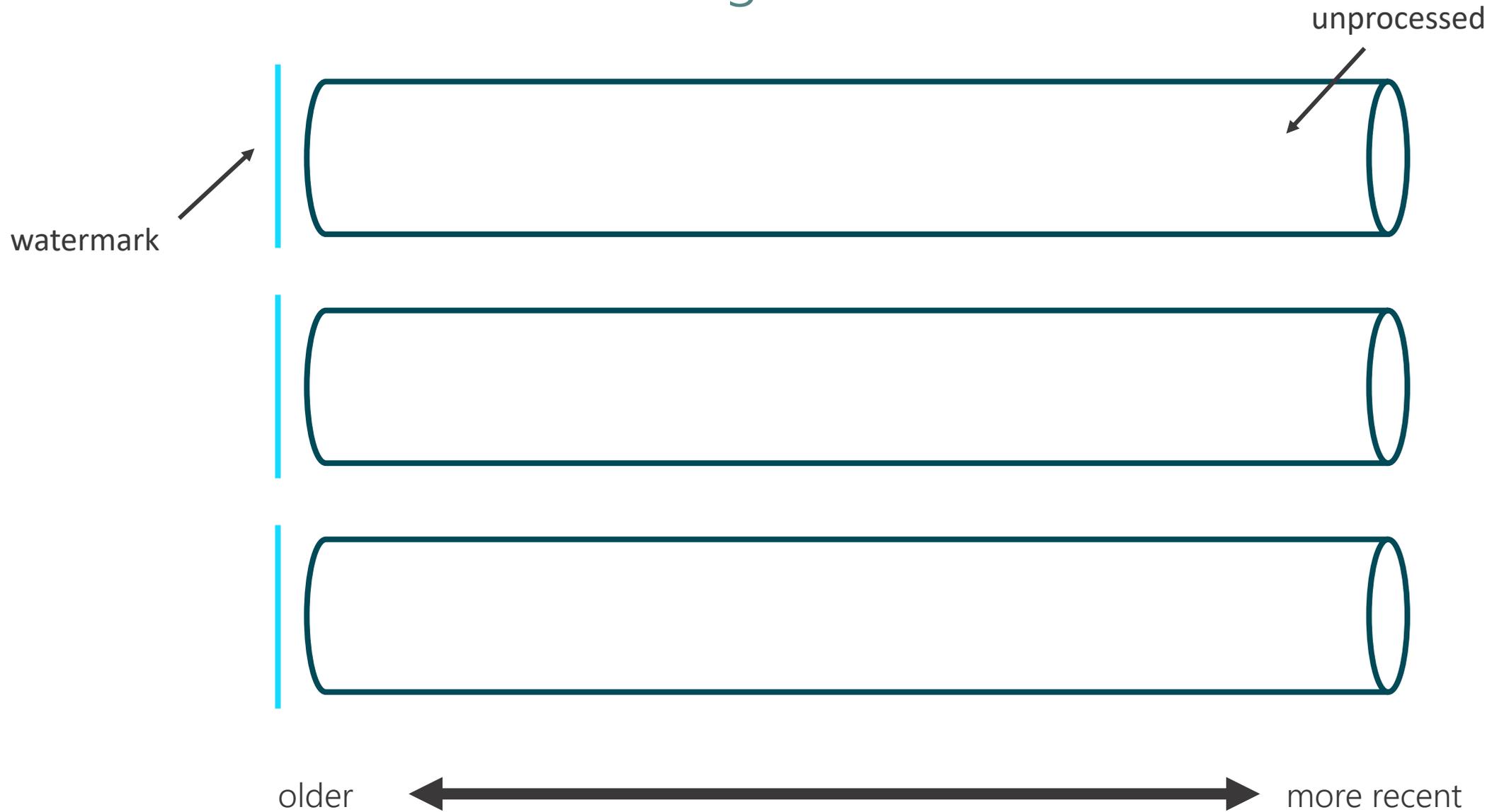Latency SLAs

Completeness and Latency is a tradeoff

## Batch Processing
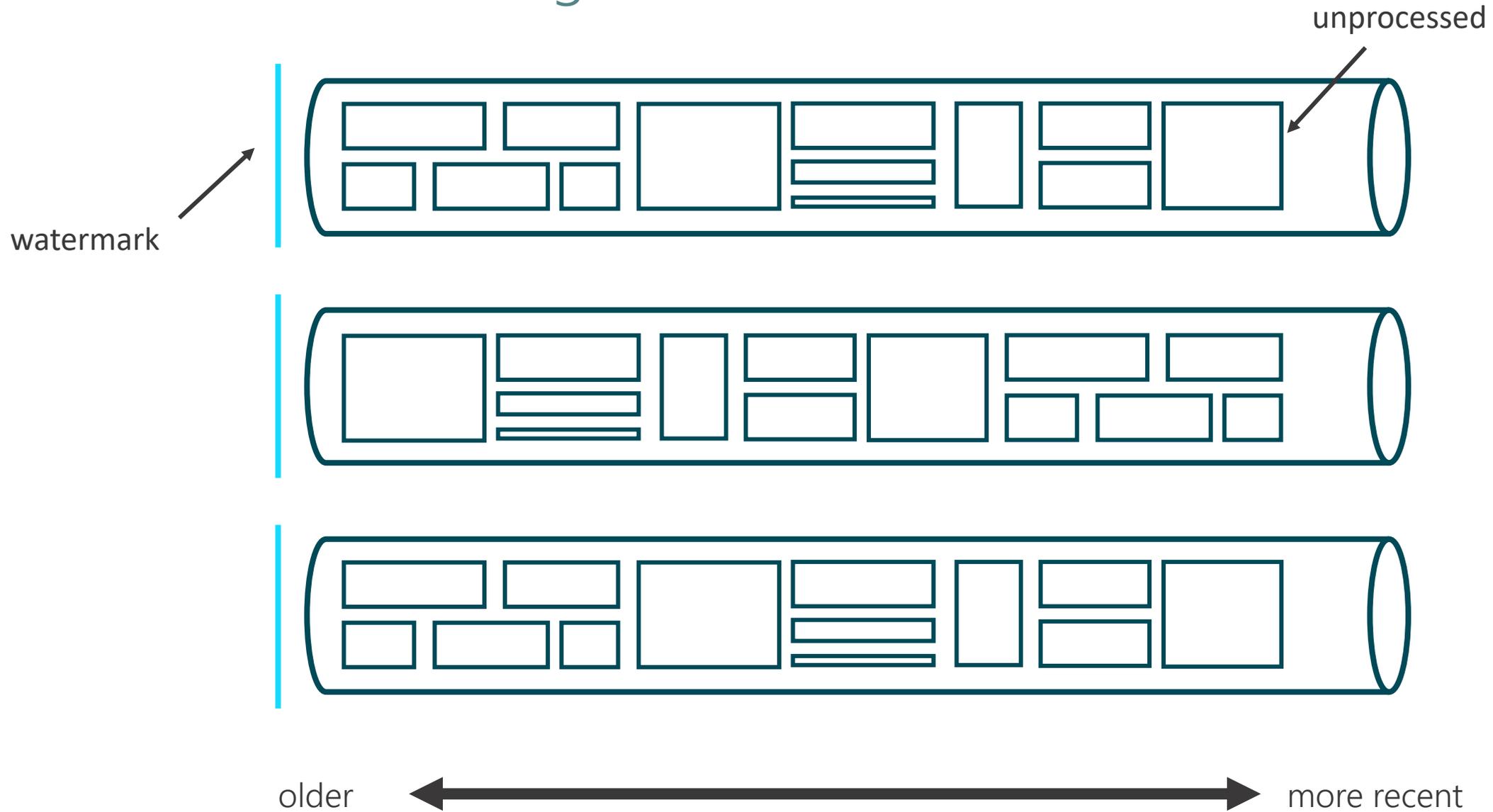
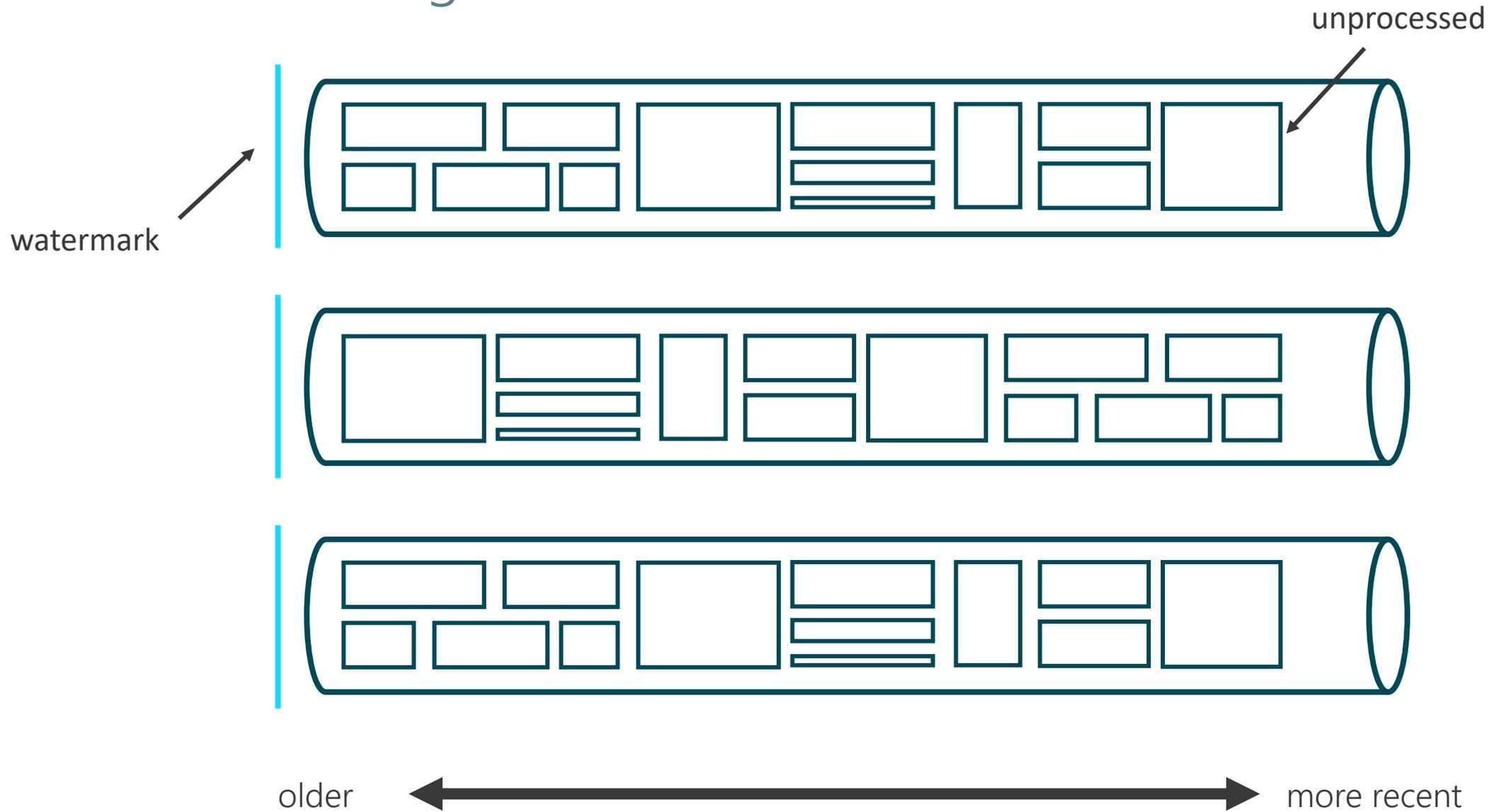Data is as complete as it gets within the job

No Low Latency SLAs

# Stream Real-time Processing

unprocessed

watermark

older ⟵⟶ more recent

|

# Stream Re-Processing

unprocessed

watermark

older ←————————→ more recent

# Batch Processing

unprocessed

watermark

older ← → more recent

# Batch vs. Stream Processing

## Continuous Streaming

Watermarks to model Completeness/Latency tradeoff

Incremental results & Proc.-Time Timers

In-receive-order ingestion with low parallelism

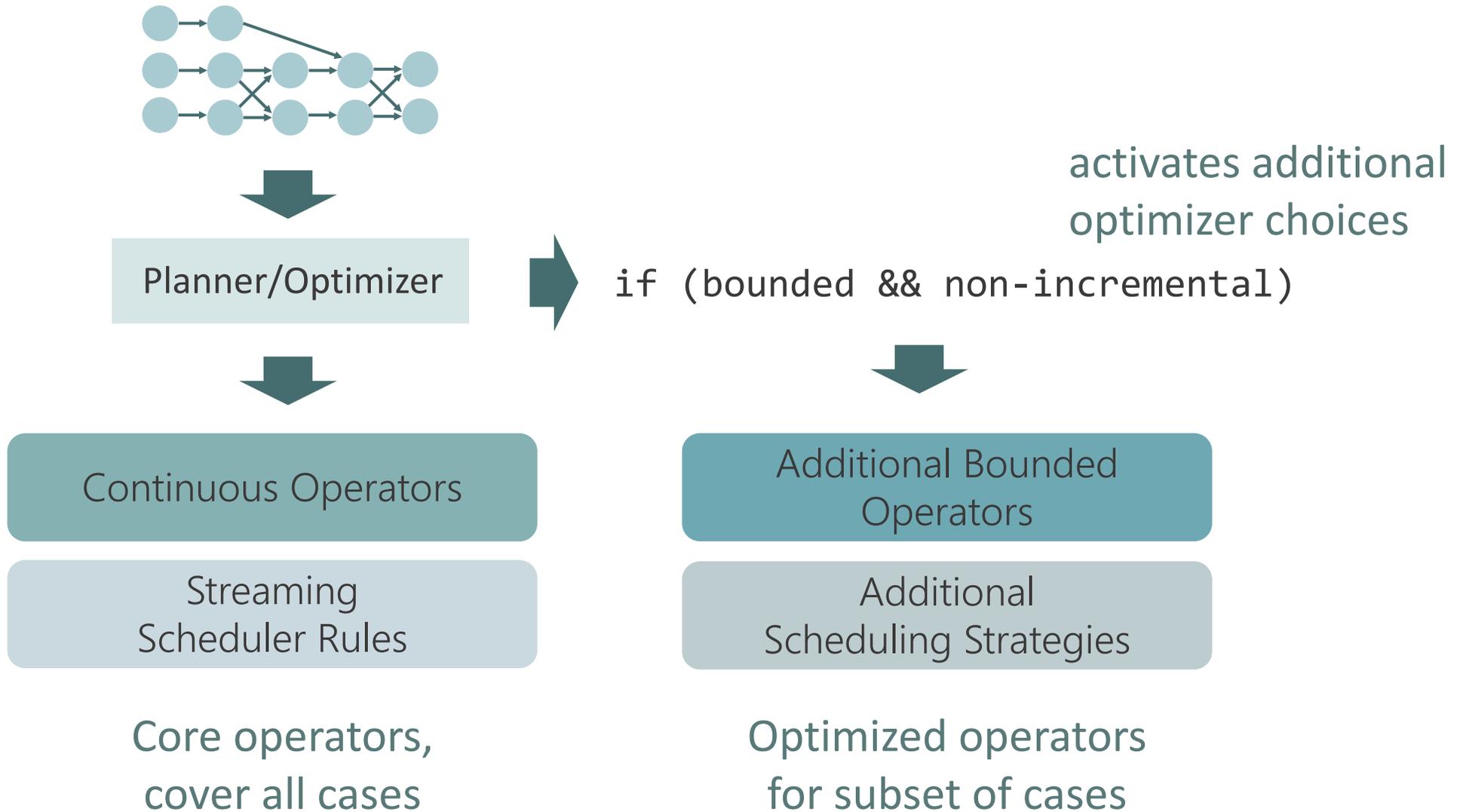## Batch Processing

No Watermarks

Results at end-of-program only
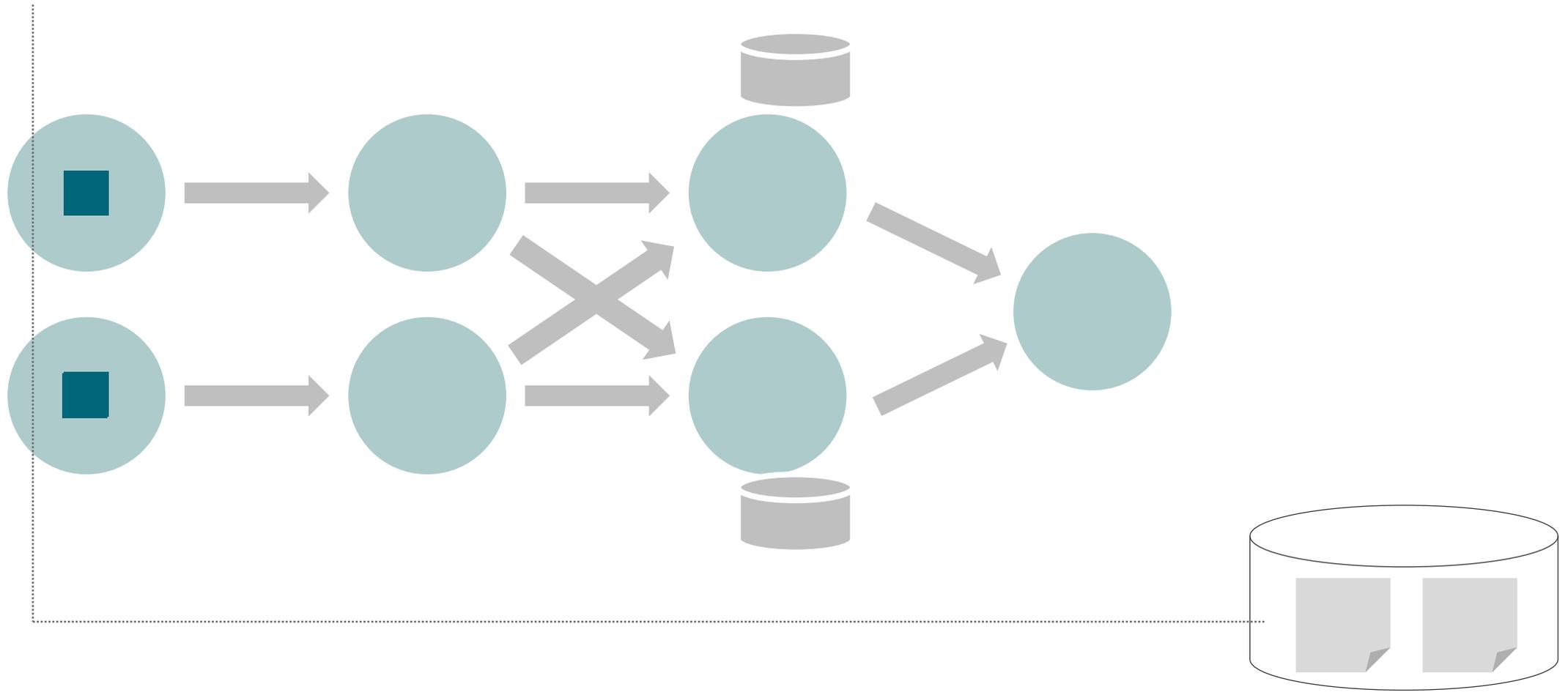
Massively parallel out-of-order ingestion

# Fast Batch Runtime in a Stream Processor

# Evolution through the entire stack!
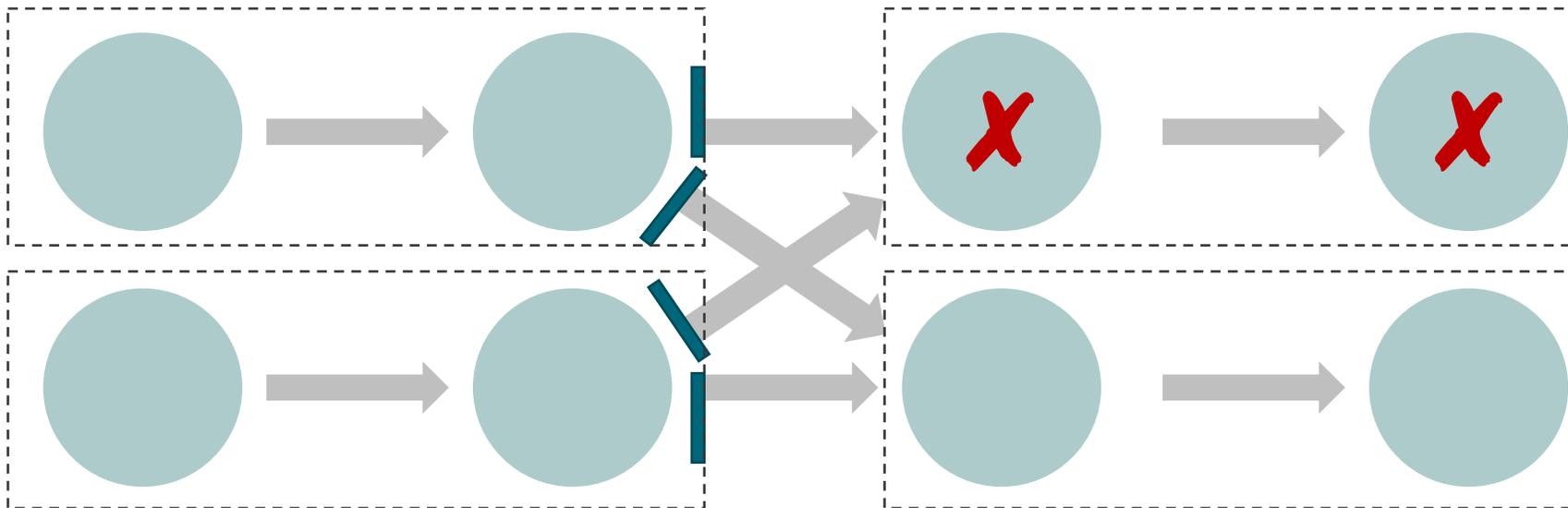
# Exploiting the Batch Special Case

activates additional optimizer choices

**Planner/Optimizer**

`if (bounded && non-incremental)`

**Continuous Operators**

**Additional Bounded Operators**

**Streaming Scheduler Rules**

**Additional Scheduling Strategies**

Core operators, cover all cases

Optimized operators for subset of cases

# Fault tolerance without writing intermediate streams to Brokers or DFS
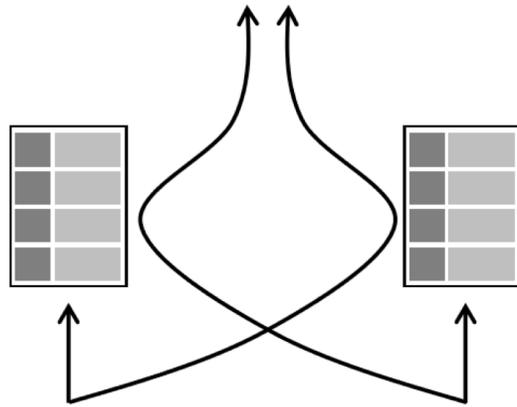
# Scheduling Strategies

- Build pipelined regions
  - Incremental results: everything pipelines
  - Non-incremental results: break pipelines once in a while

- Recovery: Restart the pipelined region from latest checkpoint (or beginning)
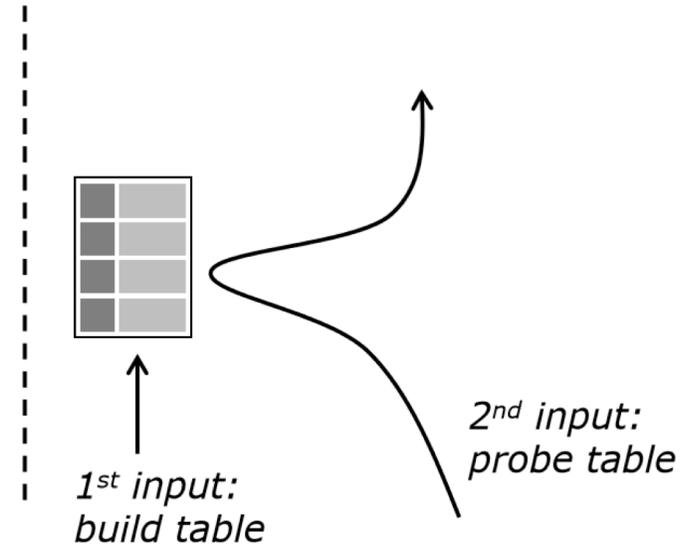  - replay input since checkpoint or beginning

# Streaming versus Batch Join



both inputs:
- build one table
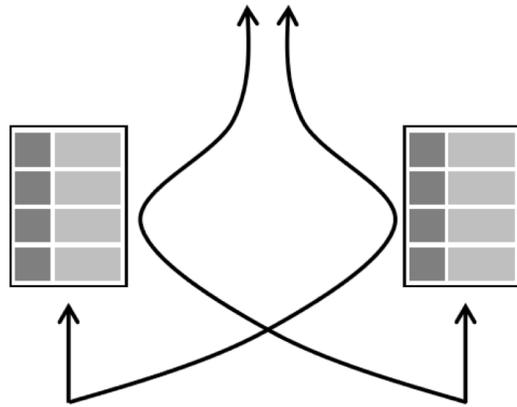- probe other table

**Continuous Streaming Join**

1st input:
build table

2nd input:
probe table

**Batch Hash Join**

|

# Streaming versus Batch Join

2x RocksDB
LSM-Trees

1x Hybrid Hash Join

bounded/
unbounded

only on
bounded data

incremental
results

both inputs:
- build one table
- probe other table

1st input:
build table

2nd input:
probe table

batch results

no checkpoints

**Continuous Streaming Join**

**Batch Hash Join**

more general

order-of-magnitude faster
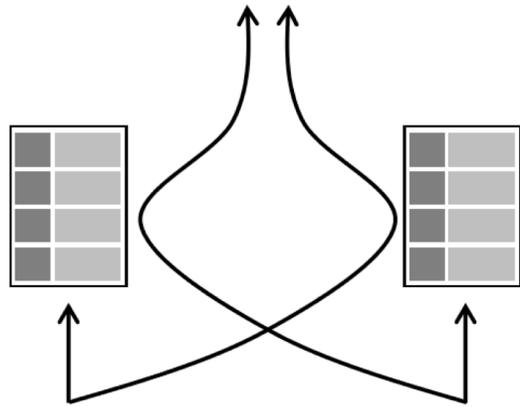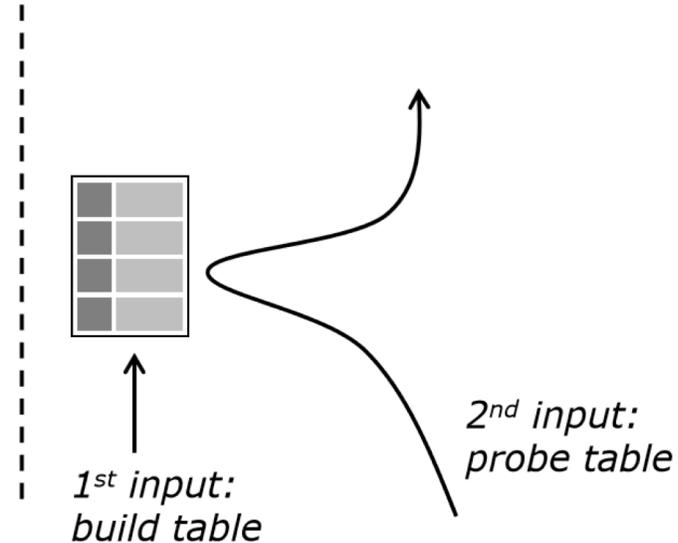
# Streaming versus Batch Join



push-based
(latency/checkpoints)

both inputs:
- build one table
- probe other table

Continuous Streaming Join

1st input:
build table

2nd input:
probe table

pull-based
(data flow control)

Batch Hash Join

more general

order-of-magnitude faster

# Push-based and Pull-based Operators

**Push Operators**

**Pull Operators**

pull()          pull()

accept data from any input immediately
(like actor messages)

pull data from one input at a time
(like reading streams)

minimize latency
supports checkpoint alignment

control over data flow,
high-latency, breaks checkpoints

# Flink 1.9 - Selectable Push-based Operators

select()                    select()

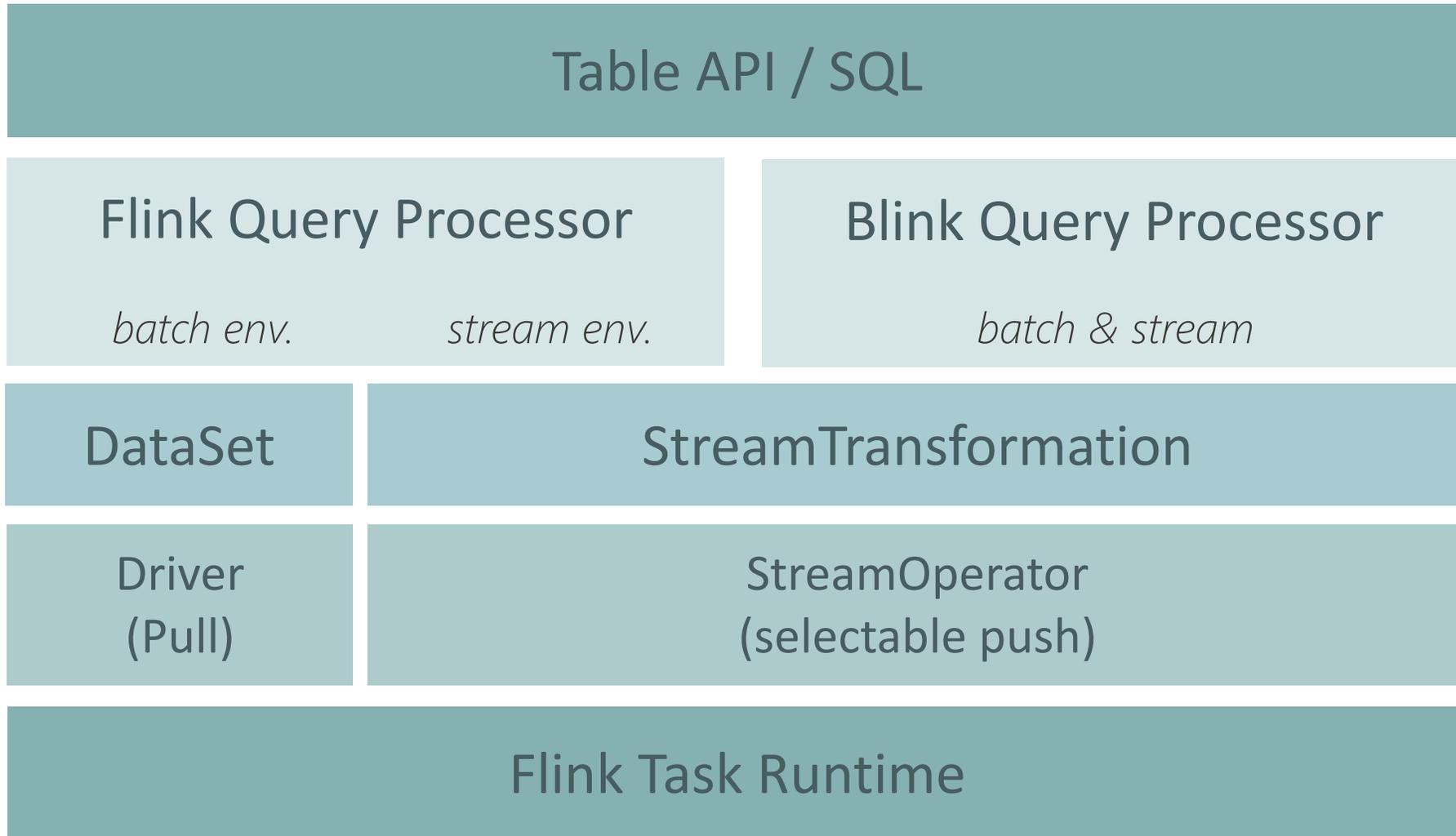similar to non-blocking-I/O model

Java NIO, Linux Epoll, or Select

subscribe to inputs (select)
and receive pushed events

➔ Operators control data flow by selecting active data paths
➔ Among active data paths, fully asynchronous data flow
   driven by network, data sources (and timers)

The State of the
Batch & Streaming Interplay in Flink

# Table API / SQL in Flink 1.9

| Table API / SQL | |
| :--- | :--- |
| **Flink Query Processor**<br><br>*batch env.*      *stream env.* | **Blink Query Processor**<br><br>*batch & stream* |

| DataSet | StreamTransformation |
| :--- | :--- |
| Driver<br>(Pull) | StreamOperator<br>(selectable push) |

| Flink Task Runtime |
| :--- |

# DataStream API

- DataStream is already supporting Bounded and Unbounded Streams

- Not exploiting batch optimizations so far
  - Bounded batch-style execution still faster on DataSet API

- After Flink 1.10:
  - Introduce `BoundedDataStream` and `non-incremental mode` to exploit optimizations for bounded data
  - Watermarks "jump" from **-∞** to **+∞** at end of program
  - Processing time timers deactivated or deferred (end of key)
  - ´Add same operators back batch-style SQL execution also for DataStream.

# Wrapping Up

# What else is new in Flink 1.9

Python Table API

Hive support

Analytics over Checkpoints/Savepoints

Preview of new Blink SQL Engine

Atomic stop-with-savepoint

...and lot's more

# What else is the community working on?

Cross-Batch-Streaming Machine Learning

Full support of Blink SQL Engine and TPC-DS coverage

New Source and Connector API

Interactive multi-job programs

Python Table UDFs

Unaligned Checkpoints

DDL and Clients for Streaming SQL

a big documentation overhaul

...and lot's more

# Learn more about Flink



Learn from the original creators of Apache Flink®

Developer Training

Barcelona, ES - Oct 21

Operations Training

Barcelona, ES - Oct 24

Register

Organized by ververica

# Or another conference?



The Apache Flink® Conference

Berlin | October 7-9, 2019

Use **FFEU19-DataCouncil** for 20% off

flink-forward.org

FLINK FORWARD

Organized by ververica

#flinkforward

# Thank you!

If you liked this, engage with the
Apache Flink® community

- Try Flink and help us improve it
- Contribute docs, code, tutorials
- Share your use cases and ideas
- Join a Flink Meetup
- Come to Flink Forward (`https://www.flink-forward.org/`)

@twalthr                    @ApacheFlink                    @VervericaData

`https://flink.apache.org/`