

DATA516/CSED516

Scalable Data Systems and Algorithms

Lecture 8

Stream Processing and Review

Final Projects

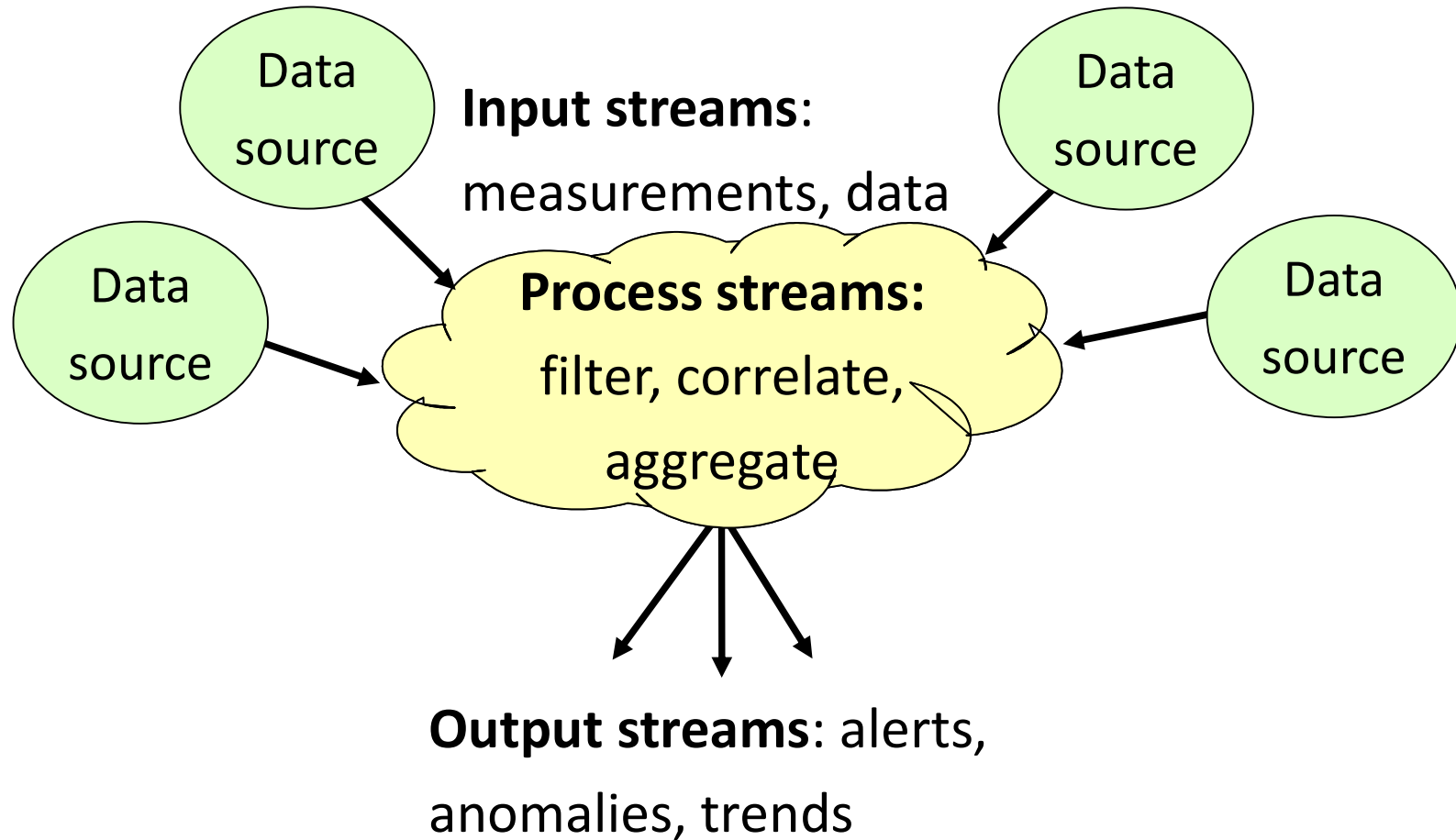
- **Deadlines**
 - HW4 Due: **Friday, December 2nd in Gitlab (+ 2 late days)**
 - Project milestone: **Friday, Nov 25th in Gitlab**
 - Project presentations: **Tuesday, Nov 29th + Dec 6th in class**
 - Final project reports: **Tuesday, December 9th in Gitlab**
 - **NO EXTENSIONS ON PROJECT DATES**
- **Project presentations (5 min)**
 - Looking for more volunteers on for Nov 29th (add to spreadsheet)
 - TA's will send out schedule and Google Drive later this week
- **Final reports (make sure to include your name)**
 - 4-5 pages in conference paper format and style
 - Suggested outline on the course website.

Stream Processing

Batch vs Stream Processing

- **Batch Processing (Databases)**
 - Data is present before queries are issued
 - Application rely on lots of stored information
 - Bounded Dataset
 - Finite Querying Time
- **Stream Processing**
 - Data is ingested and processed as it comes in
 - Applications rely on recent data and stored data
 - Unbounded Dataset
 - “Queries” can run for months/years/decades...

Stream Processing: Early Days



Application Domains

- Network monitoring
 - Intrusion, fraud, anomaly detection, click streams
- Financial services
 - Market feed processing, ticker failure detection
- Sensor-based environment monitoring
 - Weather conditions, air quality, car traffic
 - Civil engineering, military applications, etc.
- Medical applications
 - Patient monitoring, equipment tracking
- Near real-time data analytics

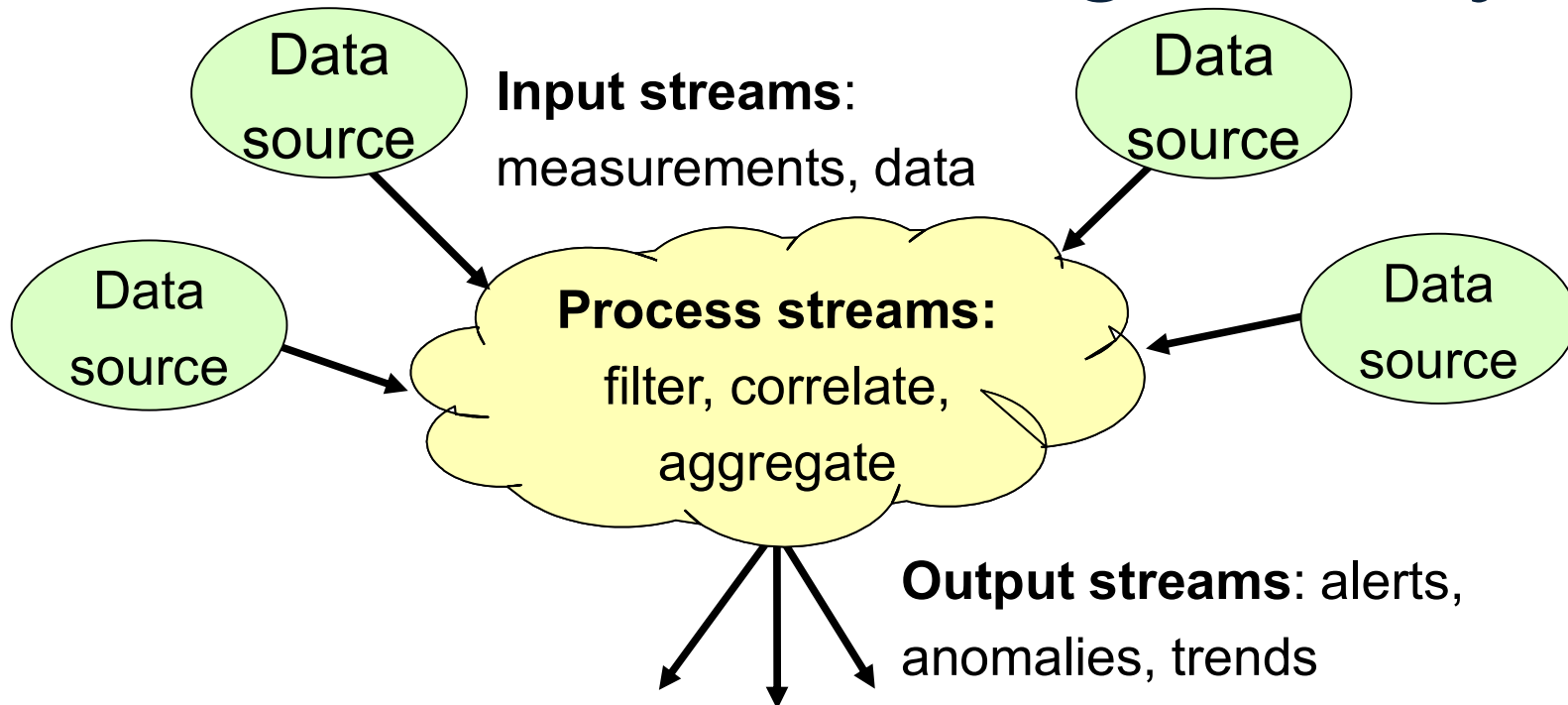
Requirements

- **Input data is pushed continuously**
 - Traditional DBMSs not designed for continuous loading or inserting of individual data items
 - “DBMS-active, human passive” model
- **Users want to execute continuous queries**
 - Traditional DBMSs have no direct support for such queries. Can use triggers, but triggers do not scale
- **Low-latency processing**
 - Need to see results in near real-time
 - Data is possibly high-volume and high-rate

Other Requirements

- Distribution
- Load management and load shedding
- Approximate processing, approximate answers
- Fault-tolerance

Stream Processing: Today



Application domains: IoT, Web analytics, application telemetry, finance, healthcare

Stream processing engines: Kafka, Heron, Trill, StreamInsight, Spark Streaming, Beam, Flink, ...

Same But Different Stream Processing

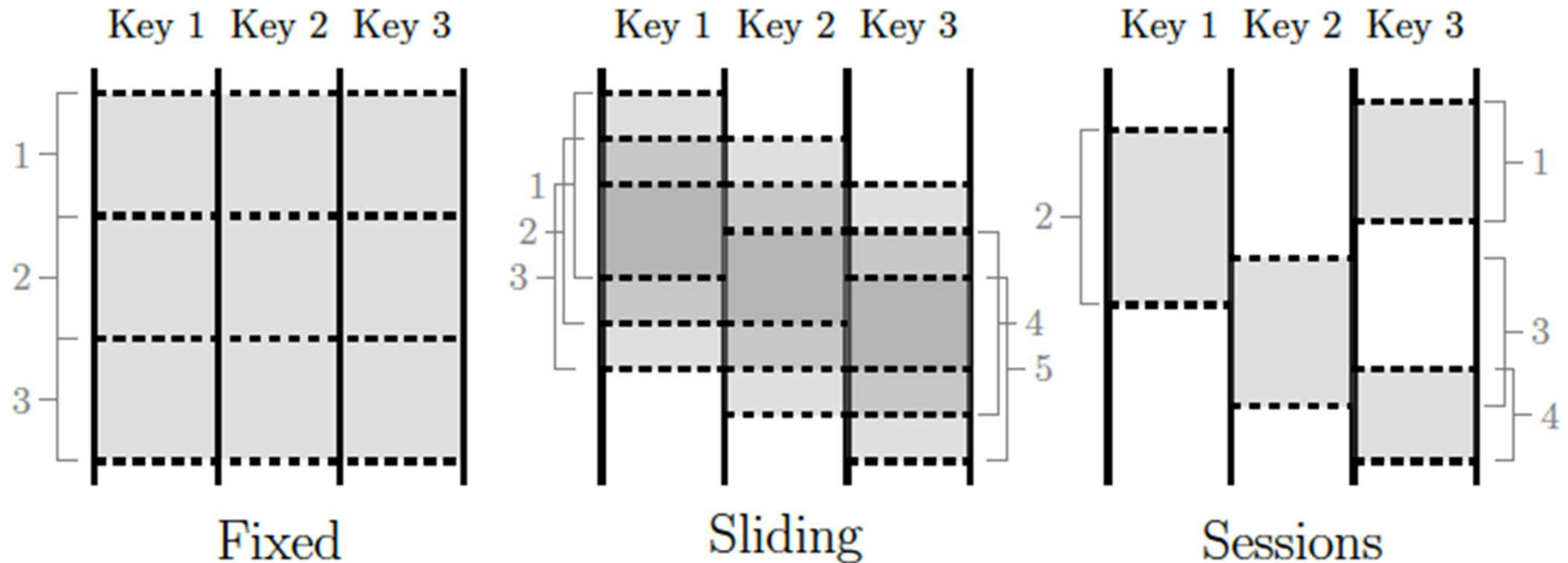
- Today, stream processing fundamentally same
 - Unbounded streams of tuples, timestamps, windows, ...
- But recent systems have different emphases
 - Single programming model for batch and streaming
 - Parallel, shared-nothing stream processing
 - APIs in Python, Java, etc
 - Seamless support for user-defined functions

Streaming Concepts

Types of Windows

- **Fixed (Tumbling) windows: Static window size**
 - Hourly windows or daily windows
- **Sliding windows: Overlapping static windows**
 - Defined by a window size and a slide interval
 - Hourly window sliding by one min
- **Sessions: Data-dependent windows**
 - Group data by key
 - For each key, a window is a burst of data in the stream
 - Window ends when a timeout occurs

Types of Window Illustration

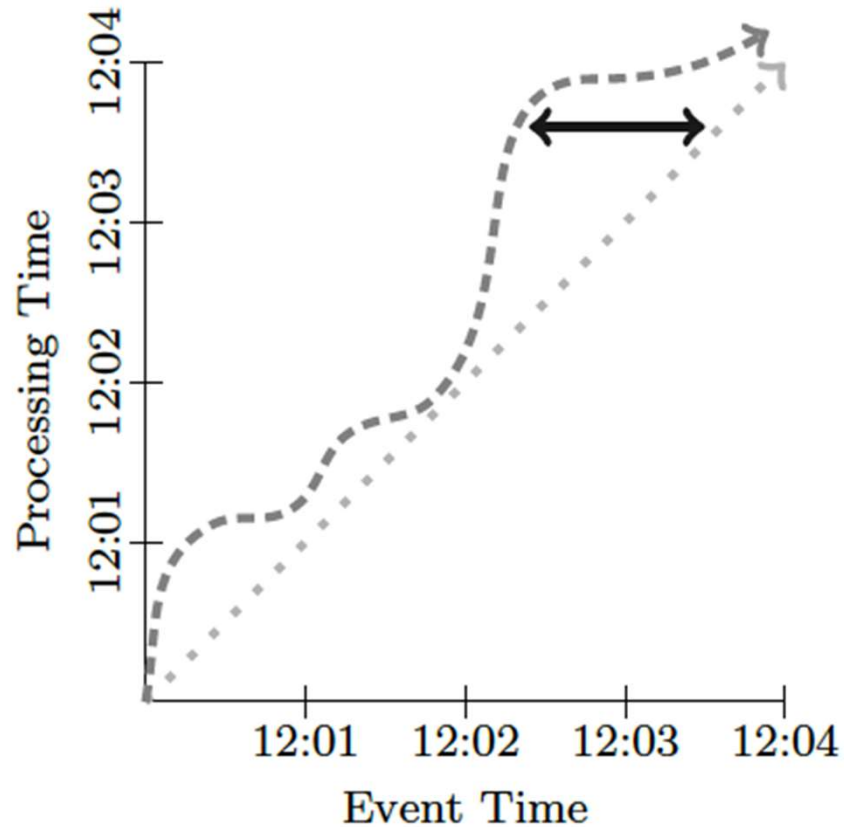


From: Akidau et. al. [The dataflow model: a practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing](#). VLDB'15.




Two Notions of Time in Streams

- **Event time**: Time when event occurred
- **Processing time**: Time when event observed
- **Time domain skew**
 - Difference between event and processing times
 - No skew: Process events immediately when they happen
- **Watermark**
 - Heuristic, lower bound on event time processed
 - Semantics: Future tuples should have higher timestamps
 - But, sometimes tuples can be late compared with watermark

Time Domain Skew Illustrated



From: Akidau et. al. [The dataflow model:...](#)
VLDB'15.

Actual watermark: 
Ideal watermark: 
Event Time Skew: 

Watermark Challenges

- Heuristic based
- If **too fast**, then data may be late
 - Tuple with event-time 9 min may arrive after watermark for event-time 10 min was emitted
- If **too slow**, may cause processing latencies
 - If we wait for a watermark before processing data such as in a window aggregation

Stream Processing Algorithms

Constraints

- Need to process elements as they arrive
- Can only use a small amount of memory
- No time to read from or write to disk

- Often, we will use *approximate* algorithms

General Sampling Approach

- To select a fraction a/b of stream elements
- One attribute in stream is key
 - The user ID in the next example
 - But it could be the search query or another attribute
- Hash key value into b buckets
- Retain all stream items that hash into first a buckets

Sampling from Streams

- Goal: Collect a representative sample of stream data
- Example:
 - Search engine receives a stream of queries
 - “What fraction of the typical user’s queries were repeated over the past month?”
 - Wish to store only 1/10th of the stream elements
- Challenge:
 - If we pick 1/10th of all queries, hard to reason about complex user behavior such as fraction of repeated queries
 - Solution: pick 1/10th of users and keep all queries for those users
 - How? Hash user IDs into 10 buckets. If a user hashes into bucket 0, keep the corresponding query

Fixed Sample Size

- What if we want to sample N items
 - As opposed to $N\%$ of all the items?
- Approach: Hash keys into B buckets
 - Make B a large number
- Every time get to $> N$ items, drop the last bucket and all values that previously hashed into that bucket

Stream Selection

- **Goal: Apply a filter to a stream**
 - If a tuple meets the selection condition, keep it
 - Otherwise, drop it
- **Challenge: Some filter predicates are expensive to compute**
 - Example: Look up email address to decide if spam
- **Solution: Bloom filters**

Bloom Filters

0	0	0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---

S: Set of key values that pass the filter

For each value v in S , compute $\text{hash}(v)$, set corresponding bit to 1

0	1	0	0	1	0	1	1	0	0
---	---	---	---	---	---	---	---	---	---

For each tuple in the stream with key value w

Compute $\text{hash}(w)$

If corresponding bit is 1, then tuple passes the filter

Bloom Filters

0	0	0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---

S: Set of key values that pass the filter

For each value v in S , compute $\text{hash}(v)$, set corresponding bit to 1

0	1	0	0	1	0	1	1	0	0
---	---	---	---	---	---	---	---	---	---

For each tuple in the stream with key value w

Compute $\text{hash}(w)$

If corresponding bit is 1, then tuple passes the filter

Not Perfect?

Bloom Filters

0	0	0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---

S: Set of key values that pass the filter

For each value v in S , compute $\text{hash}(v)$, set corresponding bit to 1

0	1	0	0	1	0	1	1	0	0
---	---	---	---	---	---	---	---	---	---

For each tuple in the stream with key value w

Compute $\text{hash}(w)$

If corresponding bit is 1, then tuple passes the filter


Improvement: Use K hash functions instead of one

For each key value, compute K hashes and check K bits

Counting Distinct Elements

- Goal: Compute the number of distinct values of an attribute in a stream
- Example: Count number of distinct visitors to website
- Challenge: What if too many distinct elements to hold in memory?
- Approach: Flajolet-Martin Algorithm

Flajolet-Martin Algorithm

- Tuple in stream with value w
- Compute $\text{hash}(w) \rightarrow 001010101011000$

- R is max tail length seen so far
- Estimate of distinct elements: 2^R
 - Divided by constant factor ~ 0.77351
- Extend to many hash functions
 - Take median of group averages

Materialize

Building with Real Data

- Sacrificing Speed
- Forgoing Features
- Compromising Cost

Materialize

Materialize is a streaming database

- Written in Rust
- Maintains results of a SQL query (a materialized view) in-memory
- Provides correct answers even as the underlying data changes

Materialized View

Traditional approaches to data:

1. Collect Data
2. Write to DB
3. Query the DB

Materialized View

Traditional approaches to data:

1. Collect Data
2. Write to DB
3. Query the DB

Expensive Operations at (3) Querying Time
(Joins, Group By, etc) increases Latency

Materialized View

Materialized Views are precomputed query results whose output is stored for fast usage

In a streaming setting, these results need to be refreshed and updated to include recent data

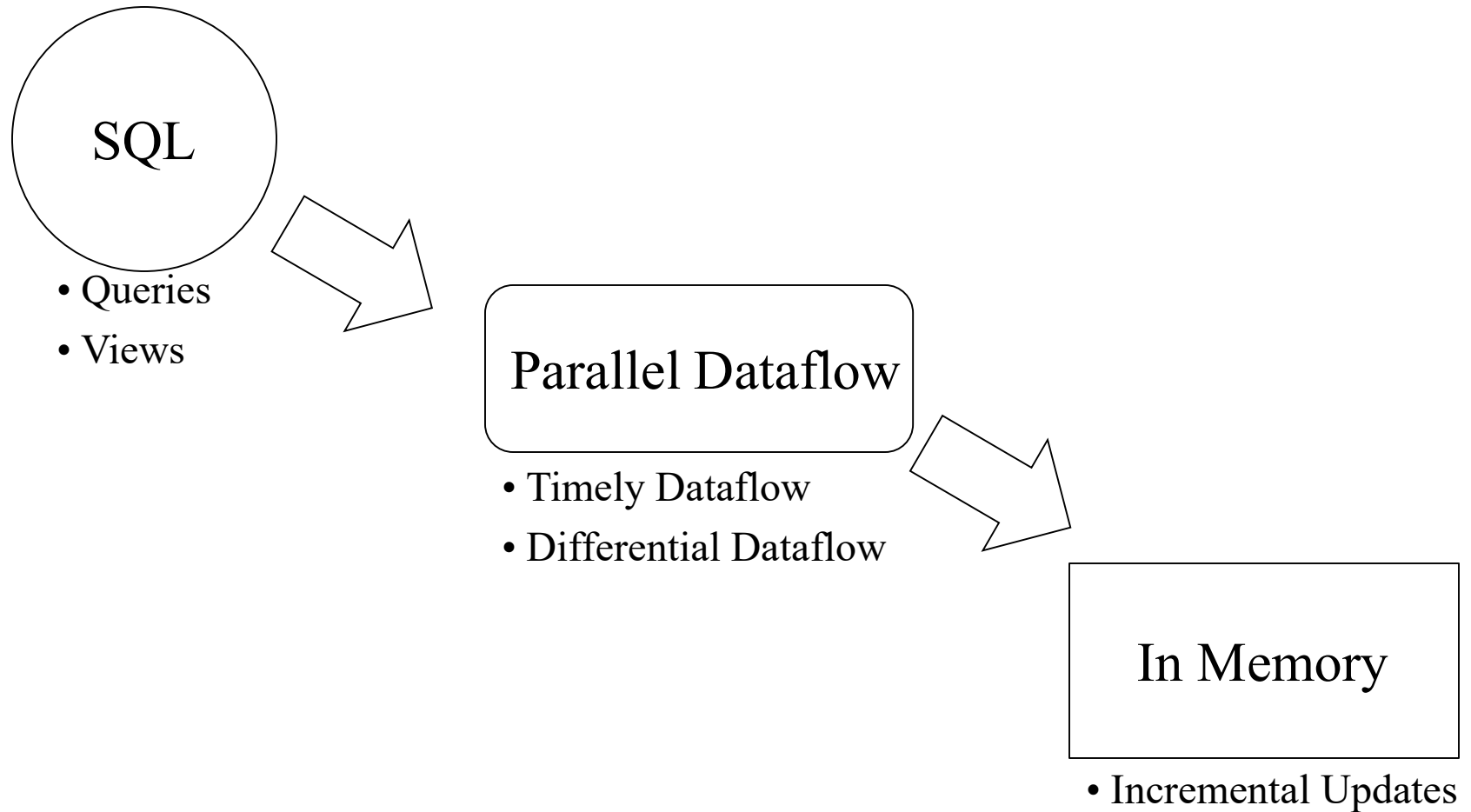
Materialized View

Materialized Views are precomputed query results whose output is stored for fast usage

In a streaming setting, these results need to be refreshed and updated to include recent data

Materialize allows definition of materialized views on incoming data and incrementally updates the views

How does Materialize Work



Dataflow (High level)

When a query/view/materialized view is created in **Materialize** it is translated to a dataflow.

Dataflow is transformation topology representing what/how the query output should be calculated

Materialize

Benefits

- Standard SQL
- Interoperability
- Complex Join Support
 - FULL, OUTER, CROSS, etc

Materialize

Syntax for Creating a Source

```
CREATE SOURCE mz_source  
FROM POSTGRES CONNECTION pg_connection  
(PUBLICATION 'mz_source')  
FOR ALL TABLES  
WITH (SIZE = '3xsmall');
```

•<https://materialize.com/docs/sql/create-source/>

Materialize

Syntax for Creating a Materialized View

```
CREATE MATERIALIZED VIEW winning_bids AS  
  
SELECT auction_id,  
        bid_id,  
        item,  
        amount  
FROM highest_bid_per_auction  
WHERE end_time < mz_now();
```

•<https://materialize.com/docs/sql/create-materialized-view/>

Materialize

Benefits

- High Performance via Dataflow
 - Work is proportional to new data (not total volume)
- Correctness over Eventual Consistency
 - Via logical update timestamps
- “Efficient Resource Consumption”

Materialize

Downsides

- If results don't need to be maintained, just use a traditional data processor
- If computations don't share states, just calculate from scratch with fastest executor

Conclusion

- Stream processing was and still is an active research area
- It is now a key component of big data solutions in industry
- Many algorithms specialized for stream processing
- Today's systems and techniques build on past work in database community

References

- Abadi et. al. [Aurora: a new model and architecture for data stream management](#). The VLDB Journal 12, 2 (2003).
- Akidau et. al. [The dataflow model: a practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing](#). VLDB'15.
- Jure Leskovec, Anand Rajaraman, and Jeffrey D. Ullman. [Mining of Massive Datasets](#). Chapter 4 (Sections 4.1 - 4.4).
- Brandon Hayes. CSED516 Guest Lecture
- Jonathan Leang. CSE344
- <https://materialize.com>

Course Review

Topics

- Relational Model
- Query Execution/Optimization
- MapReduce/Spark
- Parallel Query Evaluation
- Graphs
- Column Stores
- Streaming

Topics

- Relational Model
- Query Execution/Optimization
 - Relational Algebra
 - Data Independence
 - Logical/Physical Query Plans
 - Optimization
 - Cardinality Estimation

Topics

- MapReduce/Spark
- Parallel Query Evaluation
 - MapReduce
 - Spark
 - Parallel Query Plans
 - Hashing, Partitioning, Shuffling, etc

Topics

- Graphs
 - Recursions
 - SQL Limitations
- Column Stores
 - Design Choices
 - Efficiency
- Streaming

Tools

- AWS Redshift, S3, EC2
- ~Databricks
- Snowflake
- Docker
- Souffle
- Vertica
- Materialize