Data Integration and Large Scale Analysis

Slides credit: Matthias Boehm - Shafaq Siddiqi

11- Stream Processing





Lucas Iacono. PhD. - 2024

Part B

Large-Scale Data Management & Analysis

- LU4. Large-Scale Data Analysis
 - Distributed, Data-Parallel
 Computation [Dec 20]
 - Distributed Stream Processing
 [Jan 10]
 - Distributed Machine Learning Systems [Jan 17]

Agenda

- Data Stream Processing
- Distributed Stream Processing
- Data Stream Mining

Distributed Stream Processing

Data Stream Processing



Ubiquitous Data Streams

Event and message streams (e.g., clicks, twitter, IoT, Sensor Networks, Machines)

Q -	Characteristics	
	Event Streams	Message Streams
	 Time-oriented. Immutable Real-time processing 	 System communication Asynchronous Structured

Stream Processing Architecture

- Infinite input streams, often with windows semantics
- Continuous queries



Stream Processing Engines



Use Cases

- Monitoring and alerting (notifications on events / patterns)
- **Real-time reporting** (aggregate statistics for dashboards)
- **Real-time ETL** and event-driven data updates
- Real-time decision making (fraud detection)
- Data stream mining (summary statistics w/ limited memory)

Data Stream

- Unbounded stream of data tuples S = (s1,s2, ...) with si=(ti,di)
- S={(t1,22.5),(t2,22.7),(t3,22.8),...}

Real-time Latency Requirements

- **Real-time:** guaranteed task completion by a given deadline (30 fps)
- Near Real-time: few milliseconds to seconds
- In practice, used with much weaker meaning

Challenges in Real-Time Systems:

- Resource Constraints
- Latency
- Concurrency

Data Stream Processing: History of Stream Processing Systems

2000s

- Data stream management systems: STREAM(Stanford'01), Aurora (Brown/MIT/Brandeis'02), TelegraphCQ (Berkeley'03) but mostly unsuccessful in industry/practice
- Message-oriented middleware and Enterprise Application Integration (EAI): IBM Message Broker, SAP eXchange Infra

Academic DSMS	MOM (Message-oriented middlewares)	EAI (Enterprise Application Integration)	
Real-time data stream processing	Asynchronous message communication between systems	Workflow orchestration for enterprise systems	
SQL-like queries on streams	Sending simple or complex messages	Integration of complex applications	
Time windows and event-based operations	Event-based message delivery	Event-driven task coordination	
Academic prototypes	Industrial implementations (Kafka)	Enterprise suites IBM MQ	

Data Stream Processing: History of Stream Processing Systems

2010s

- **Distributed stream processing engines**, and "unified" batch/stream processing
- **Proprietary systems:** Google Cloud Dataflow, MS StreamInsight / Azure Stream Analytics, IBM InfoSphere Streams / Streaming Analytics, AWS Kinesis
- Open-source systems: Apache Spark Streaming (Databricks), Apache Flink (Data Artisans), Apache Kafka (Confluent), Apache Storm









System Architecture - Native Streaming



System Architecture - Native Streaming

Basic System Architecture

- Data flow graphs (potentially w/ multiple consumers)
- Nodes asynchronous ops (w/ state) (e.g., separate threads)
- **Edges** data dependencies (tuple/message streams)
- **Push model** data production controlled by source

Operator Model

- Read from input queue
- Write to potentially many output queues



Multi-Query Optimization

- Given set of continuous queries compile minimal DAG w/o redundancy -> subexpression elimination -> avoid redundant operations and share intermediate results between queries to improve system efficiency.
- Operator and Queue Sharing

Multi-Query Optimization







Multi-Query Optimization

Operator Sharing: complex ops w/ multiple predicates for adaptive reordering

Queue Sharing: share results with multiple queries



Operator sharing: complex ops w/ multiple predicates for adaptive reordering

// Execute both queries
EXECUTE average_pressures
EXECUTE count_pressures

Same filter applied twice in Q1 & Q2 -> duplicated work.

 Operator sharing	<pre>//Apply the common filter only once filtered_pressures = FILTER "pressures" WHERE value > 1000</pre>
<pre>// Query 1: Average of pressures greater than 1 average_pressures = FILTER "pressures" WHERE va > 100 GROUP BY 1-minute_interv COMPUTE AVERAGE(value)</pre>	<pre>// Query 1: Average value of the filtered pressures average_pressures = GROUP filtered_pressures BY 1-minute_interval</pre>
<pre>// Query 2: Count of pressures greater than 100 count_pressures = FILTER "pressures" WHERE valu 1000</pre>	<pre>// Query 2: Count of the filtered pressures count_pressures = GROUP filtered_pressures BY 1-minute_interval</pre>

- Back Pressure
 - Graceful handling of overload
 w/o data loss
 - \circ Slow down sources
 - \circ $\,$ E.g. blocking queues



Back Pressure

- Graceful handling of overload
 w/o data loss
- \circ Slow down sources
- E.g. blocking queues



Self-adjusting operator scheduling Pipeline runs at rate of **slowest op**

Back Pressure

- Graceful handling of overload Ο w/o data loss
- Slow down sources 0
- E.g. blocking queues 0

Load Shedding

- Random-sampling-based load shedding Ο
- **Relevance-based** load shedding 0
- **Summary-based** load shedding (synopses) 0
- Given SLA, select queries and shedding placement that minimize error [Nesime Tatbul Ο et al: Load and satisfy constraints Shedding in a Data Stream

Α

3ms 9ms Self-adjusting operator scheduling Pipeline runs at rate of slowest op

В



Manager. VLDB

2003]

2ms

Back Pressure

- Graceful handling of overload
 w/o data loss
- Slow down sources
- E.g. blocking queues

• Load Shedding

- Random-sampling-based load shedding
- Relevance-based load shedding
- Summary-based load shedding (synopses)
- Given SLA, select queries and shedding placement that minimize error and satisfy constraints

Α

3ms

• Distributed Stream Processing

- Data flow partitioning (distribute the **query**)
- Key range partitioning (distribute the **data** stream)

Self-adjusting operator scheduling Pipeline runs at rate of slowest op

В

9ms



et al: Load

2003]

Shedding in a Data Stream Manager. VLDB

2ms

Time (Event, System, Processing)

- Event Time
 - Real time when the event/data item was created
- Ingestion
 - \circ $\,$ System time when the data item was received
- Processing Time
 - \circ $\,$ System time when the data item is processed



Event Time

Time (Event, System, Processing)

- Event Time
 - Real time when the event/data item was created
- Ingestion
 - \circ $\,$ System time when the data item was received
- Processing Time
 - \circ $\,$ System time when the data item is processed





Time (Event, System, Processing)

- Event Time
 - Real time when the event/data item was created
- Ingestion
 - \circ System time when the data item was received
- Processing Time
 - System time when the data item is processed
- In practice
 - Delayed and unordered data items
 - Use of heuristics (e.g watermarks, delays)
 - Use of more complex triggers (speculative and late results)





Durability and Consistency Guarantees

- At Most Once
 - "Send and forget", ensure data is never counted twice
 - Might cause data loss on failures

03 Message-oriented Middleware, EAI, and Replication







Durability and Consistency Guarantees

- At Most Once
 - "Send and forget", ensure data is never counted twice
 - \circ $\,$ Might cause data loss on failures $\,$
- At Least Once
 - "Store and forward" deliver the message until reception of the acknowledgements from receiver
 - Might create incorrect state (processed multiple times)

03 Message-oriented Middleware, EAI, and Replication





Durability and Consistency Guarantees

- At Most Once
 - "Send and forget", ensure data is never processed twice
 - \circ $\,$ Might cause data loss on failures $\,$
- At Least Once
 - "Store and forward" deliver the message until reception of the acknowledgements from receiver
 - Might create incorrect state (processed multiple times)
- Exactly Once
 - "Store and forward" w/ guarantees regarding state updates and sent msgs
 - Often via dedicated transaction mechanisms (hand-shaking protocols)

03 Message-oriented Middleware, EAI, and Replication





Window Semantics

• Windowing Approach

_ _ _

- Many operations like joins/aggregation **undefined over unbounded streams**
- \circ Compute operations over windows of (a) time or (b) elements counts

Window Semantics

• Windowing Approach

- Many operations like joins/aggregation **undefined over unbounded streams**
- \circ Compute operations over windows of (a) time or (b) elements counts

• Tumbling Window

- Every data item is only part of a single window
- Aka Jumping window



Window Semantics

• Windowing Approach

- Many operations like joins/aggregation undefined over unbounded streams
- \circ Compute operations over windows of (a) time or (b) elements counts

• Tumbling Window

- Every data item is only part of a single window
- Aka Jumping window
- Sliding Window
 - \circ $\,$ Time- or tuple-based sliding windows $\,$
 - Insert new and expire old data items



Basic Stream Join

_ _ _

- **Tumbling window:** use classic join methods
- **Sliding window** (symmetric for both R and S)
 - Applies to arbitrary join pred

Basic Stream Join

- **Tumbling window:** use classic join methods
- **Sliding window** (symmetric for both R and S)
 - Applies to arbitrary join pred

Example: join two streams where an event in stream R (sensor_reading) should be matched with a value in stream S (control command) in a 3-second sliding window.

For each new r in R:

- a. **Scan** window of stream S to find match tuples
- b. Insert new r into window of stream R
- c. Invalidate expired tuples in window of stream R

- Double-Pipelined Hash Join
 - Join of bounded streams (or unbounded w/ invalidation)
 - Equi join predicate, symmetric and non-blocking
 - For every incoming tuple (e.g. left): probe (right)+emit, and build (left)



- Double-Pipelined Hash Join
 - Join of bounded streams (or unbounded w/ invalidation)
 - Equi join predicate, symmetric and non-blocking
 - For every incoming tuple (e.g. left): probe (right)+emit, and build (left)



- Double-Pipelined Hash Join
 - Join of bounded streams (or unbounded w/ invalidation)
 - Equi join predicate, symmetric and non-blocking
 - For every incoming tuple (e.g. left): probe (right)+emit, and build (left)



- Double-Pipelined Hash Join
 - Join of bounded streams (or unbounded w/ invalidation)
 - Equi join predicate, symmetric and non-blocking
 - For every incoming tuple (e.g. left): probe (right)+emit, and build (left)



HR, RID Hs, sid ► RID=SID Double-Pipelined Hash Join 1,1,2 1,7 of bounded 0 Join streams (or unbounded w/ invalidation) Equi join predicate, 0 symmetric and non-blocking ab 7 1 zy incoming Ο For every (e.g. left): tuple 2 cd 1 XW probe (right)+emit, ef emit 1(efxw) 1 and build (left)

> Stream Stream R S

emit 1(abxw)

- Double-Pipelined Hash Join
 - Join of bounded streams (or unbounded w/ invalidation)
 - Equi join predicate, symmetric and non-blocking
 - For every incoming tuple (e.g. left): probe (right)+emit, and build (left)



HR, RID Hs, sid ► RID=SID Double-Pipelined Hash Join 1,1,2,7 1,7,7 of bounded 0 Join streams (or unbounded w/ invalidation) Equi join predicate, 0 symmetric and non-blocking ab 7 1 ZY incoming Ο For every (e.g. left): tuple 2 cd 1 XW emit 1(abxw) (right)+emit, probe ef 7 1 emit 7(ghvu) emit 1(efxw) vu and build (left) 7 gh emit 7(ghzy) Stream Stream R S

Distributed Stream Processing

Query-Aware Stream Partitioning

Example Use Case

- AT&T network monitoring with Gigascope (e.g., OC768 network)
- 2x40 Gbit/s traffic 112M packets/s 26 cycles/tuple on 3Ghz CPU
- Complex query sets (apps w/ ~50 queries) and massive data rates



T. Johnson et.al, Query-aware partitioning for monitoring massive network data streams. **SIGMOD 2008**

Query-Aware Stream Partitioning

T. Johnson et.al, Query-aware partitioning for monitoring massive network data streams. **SIGMOD 2008**



Baseline Query Execution Plan

Self join	↓ M _{tb=tb+1}	Query FLOW PAIRS: users generating heavy-loads during an specific period
High-level aggregation	γ ₂	Query HEAVY FLOWS: maximum queries per IP
Low-level aggregation	γ1	
Low-level filtering	σ Ι ΤCP	Query FLOWS: composed by 3 sub-queries about how many requests are made by each IP in the network (per minute)

Query-Aware Stream Partitioning

T. Johnson et.al, Query-aware partitioning for monitoring massive network data streams. **SIGMOD 2008**



Solution divide in sub-queries and distribute

Optimized Plan:

- Distributed Plan operators
- Pipeline and task parallelism
- Not always enough

Self join Self join High-level aggregation Low-level aggregation Low-level filtering σ

тср



Stream Group Partitioning

• Large-Scale Stream Processing

- Limited pipeline parallelism and task parallelism (independent subqueries)
- Combine with data-parallelism over stream groups
- Shuffle Grouping
 - Tuples are randomly distributed across consumer tasks

 \circ Good load balance

Stream Group Partitioning

• Fields Grouping

- \circ $% \ensuremath{\mathsf{Tuples}}$ Tuples partitioned by grouping attributes
- Guarantees order within keys, but load imbalance if skew _____

Worker

Worker

Worker

Source

Source

Stream

• Partial Key Grouping

- Apply "power of two choices" to streaming
- Key splitting: select among 2 candidates per key (works for all associative aggregation functions)



Example Apache Storm

- Example Topology DAG
 - Spouts: sources of streams
 - **Bolts:** UDF compute ops
 - Tasks mapped to worker processes and executors (threads)





Example Apache Storm

- Example Topology DAG
 - **Spouts:** sources of streams
 - **Bolts:** UDF compute ops
 - Tasks mapped to worker processes and executors (threads)



```
Config conf = new Config();
conf.setNumWorkers(3);
topBuilder.setSpout("Spout1", new FooS1(), 2);
topBuilder.setBolt("Bolt1", new FooB1(), 3).shuffleGrouping("Spout1");
topBuilder.setBolt("Bolt2", new FooB2(), 2).shuffleGrouping("Spout1");
topBuilder.setBolt("Bolt3", new FooB3(), 2)
.shuffleGrouping("Bolt1").shuffleGrouping("Bolt2");
StormSubmitter.submitTopology(..., topBuilder.createTopology());
```

Example Twitter Heron

- Motivation
 - Heavy use of Apache Storm at Twitter
 - Issues: debugging, performance, shared cluster resources, back pressure mechanism

STORM @TWITTER



Sanjeev Kulkarni et

Stream

Scale. SIGMOD 2015

Twitter Heron:

Processing at

al:

Example Twitter Heron

- Motivation
 - Heavy use of Apache Storm at Twitter
 - Issues: debugging, performance, shared cluster resources, back pressure mechanism
- Twitter Heron
 - API-compatible distributed streaming engine
 - **De-facto streaming engine at Twitter** since 2014

STORM @TWITTER Data per day Cluster Size # of formation of the second second



Example Twitter Heron

- Motivation
 - Heavy use of Apache Storm at Twitter
 - Issues: debugging, performance, shared cluster resources, back pressure mechanism
- Twitter Heron
 - API-compatible distributed streaming engine
 - **De-facto streaming engine at Twitter** since 2014
- Dhalion (Heron Extension)
 - Automatically reconfigure Heron topologies to meet throughput SLO
- Now back pressure implemented in Apache Storm 2.0 (May 2010)

STORM @TWITTER



Sanjeev

Stream

Scale.

al:

Kulkarni et

Twitter Heron:

Processing at

SIGMOD 2015

Discretized Stream (Batch) Computation

- Motivation
 - Fault tolerance (low overhead, fast recovery)
 - Combination w/ distributed batch analytics

Matei Zaharia et al: Discretized streams: fault-tolerant streaming computation at scale. SOSP 2013



Discretized Stream (Batch) Computation

- Motivation
 - Fault tolerance (low overhead, fast recovery)
 - Combination w/ distributed batch analytics
- Discretized Streams (DStream)
 - \circ Batching of input tuples (100ms 1s) based on ingest time.
 - Periodically run distributed jobs of stateless, deterministic tasks ->
 DStreams
 - State of all tasks materialized as RDDs, recovery via lineage



Matei Zaharia et al: Discretized streams: fault-tolerant streaming computation at scale. SOSP 2013



Unified Batch/Streaming Engines

- Apache Spark Streaming (Databricks)
 - Micro-batch computation with exactly-once guarantee
 - Back-pressure and water mark mechanisms
 - \circ Structured streaming via SQL (2.0), continuous streaming (2.3)
- Apache Flink (Data Artisans, now Alibaba)
 - \circ $% \ensuremath{\mathsf{Tuple}}$ Tuple-at-a-time with exactly-once guarantee
 - \circ $\,$ Back-pressure and water mark mechanisms $\,$
 - \circ $\,$ Batch processing viewed as special case of streaming $\,$





Unified Batch/Streaming Engines

- Google Cloud Dataflow
 - Tuple-at-a-time with exactly-once guarantee
 - MR FlumeJava MillWheel Dataflow
 - Google's fully managed batch and stream service
- Apache Beam (API+SDK from Dataflow)
 - Abstraction for Spark, Flink, Dataflow w/ common API, etc
 - Individual runners for the different runtime frameworks





Data Stream Mining

Overview Stream Mining

• Streaming Analysis Model

_ _ _

- \circ $% \ensuremath{\mathsf{Independent}}$ of actual storage model and processing system
- Unbounded stream of data item S = (s1, s2, ...)
- Evaluate function **f(S)** as aggregate over stream or window of stream
- Standing vs ad-hoc queries
- Recap: Classification of Aggregates
 - Additive aggregation functions (SUM, COUNT)
 - Semi-additive aggregation functions (MIN, MAX)
 - Additively computable aggregation functions (AVG, STDDEV, VAR)
 - Aggregation functions (e.g SORTING) -> approximations
- **Example:** Approximate # Distinct Items (e.g., KMV)

Stream Mining KMV

• Definition

Beyer et al. On synopses for distinct-value estimation under multiset operations. SIGMOD 2007



- Estimate the number of unique elements (cardinality) in a large data stream.
- KMV stores only a small, fixed-size subset of values derived from a hash function and uses that subset to make an accurate estimate.
- How it Works:
 - **Hashing:** Each stream element is passed through a hash function that generates a numeric value between **0 and 1**.
 - Store K Smallest Hashes: KMV keeps only the K smallest hash values.
 - \circ **Estimate** Cardinality \hat{N} using the formula:

$$\hat{N} \approx \frac{K-1}{h_{(K)}}$$

Stream Mining KMV

Data Stream: $\{A, B, C, A, D, B, E\}$

Hash Values:

$$h(A) = 0.15$$

$$h(B) = 0.45$$

$$h(C) = 0.22$$

$$h(D) = 0.70$$

$$h(E) = 0.10$$

Select K = 3: The 3 smallest hash values are: 0.10, 0.15, 0.22.

Estimate:

The largest hash in this set is 0.22, so:

$$\hat{N} \approx \frac{K-1}{h_{(K)}} = \frac{3-1}{0.22} \approx 9.09$$
 (estimated distinct values)

Note: The real number of distinct values is 5 ({A, B, C, D, E}), and the difference comes from the small K. Increasing K improves accuracy.

Beyer et al. On synopses for distinct-value estimation under multiset operations. SIGMOD 2007

Summary and Q&A

Summary and Q&A

• Summary and Q&A

_ __ __

- \circ Data Stream Processing
- \circ $\;$ Distributed Stream Processing
- \circ Data Stream Mining

• Next Lectures

- Distributed Machine Learning [Jan 17]
- Written Exam [Feb 07]

Vielen Dank!