# CSE 544
# Principles of Database Management Systems

Magdalena Balazinska

Winter 2009

Lecture 17 - Stream Processing

# Announcements

- Remember the final exam this Wednesday
  - During class
  - Open notes and open books

- Next week: project presentations Monday and Thursday
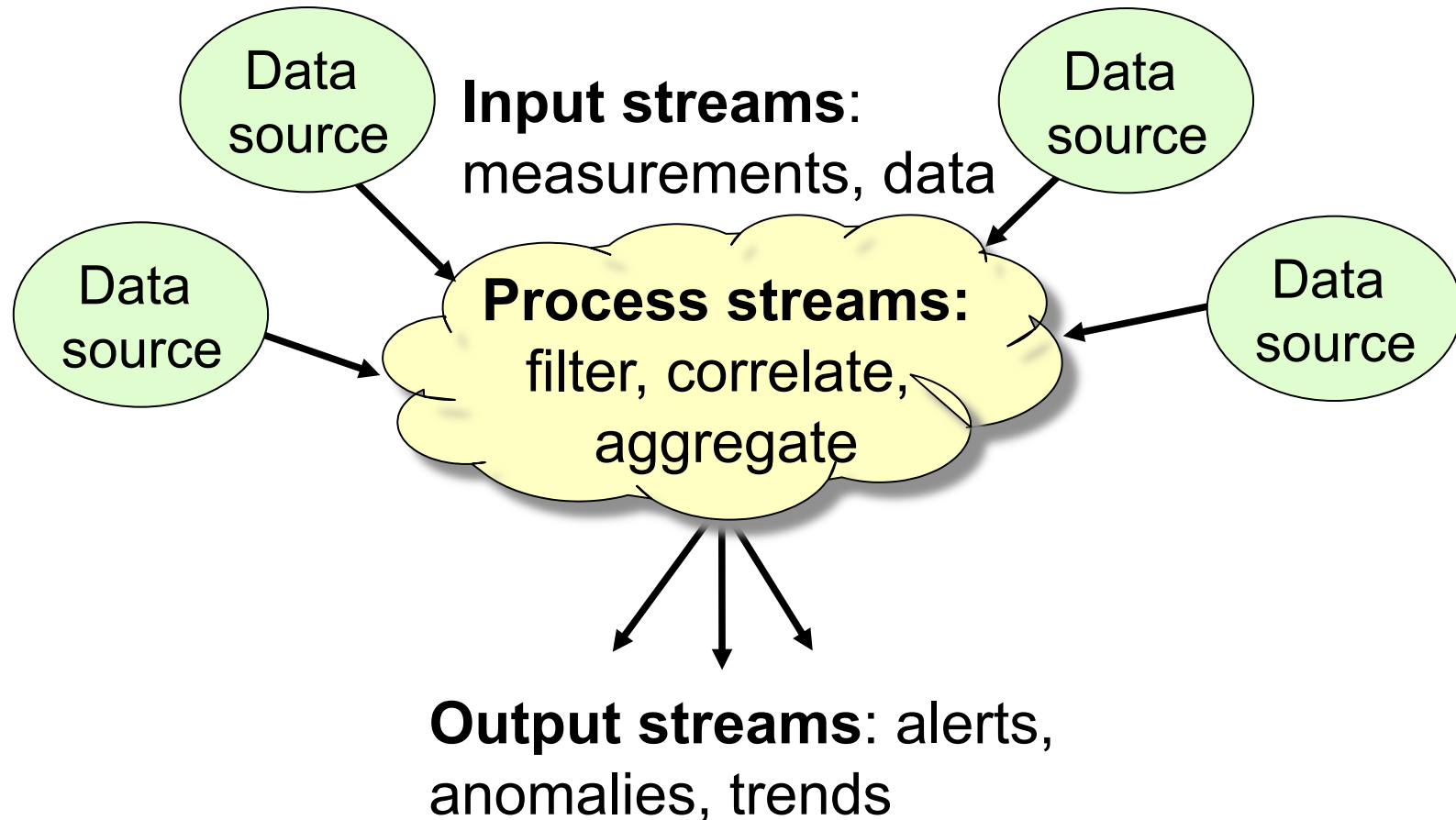
- Final reports due on Friday

# References

- **Aurora: A New Model and Architecture for Data Stream Management.** Daniel Abadi et. al. VLDB Journal. 12(2). 2003

# Outline

- **Stream processing applications**
  - Examples
  - Requirements

- **The Aurora stream processing engine**
  - Stream model and query model
  - Processing model
  - Operators
  - Query examples
  - Other features

# Stream Processing

Data
source

Data
source

Data
source

Data
source

**Input streams**:
measurements, data

**Process streams:**
filter, correlate,
aggregate

**Output streams**: alerts,
anomalies, trends

# 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 and revision processing

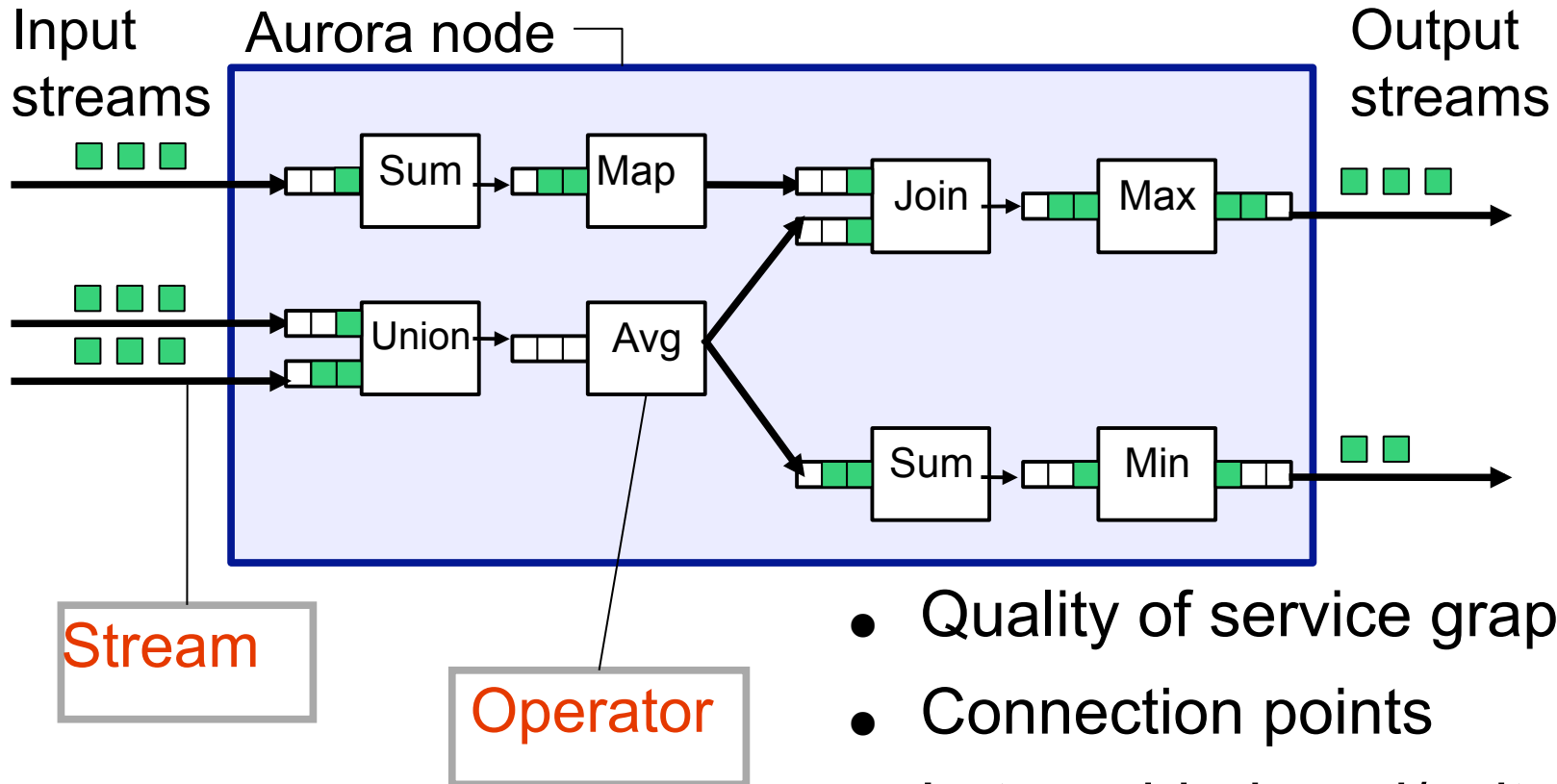- Exploiting data archives

# Outline

- **Stream processing applications**
  - Examples
  - Requirements

- **The Aurora stream processing engine**
  - Stream model and query model
  - Processing model
  - Operators
  - Query examples
  - Other features

# Stream Data Model

$$\text{Tuple}: (\underbrace{\texttt{timestamp,}}_{\text{header}} \underbrace{\texttt{v}_1\texttt{,...,v}_n}_{\text{data}})$$

- Stream: append-only sequence of tuples
- All tuples on a stream have same schema
- Timestamp is used for QoS

# Query Model



Input streams     Aurora node                         Output streams

Sum → Map → Join → Max

Union → Avg → Sum → Min

**Stream**

**Operator**

- Quality of service graphs
- Connection points
- Later added read/write ops
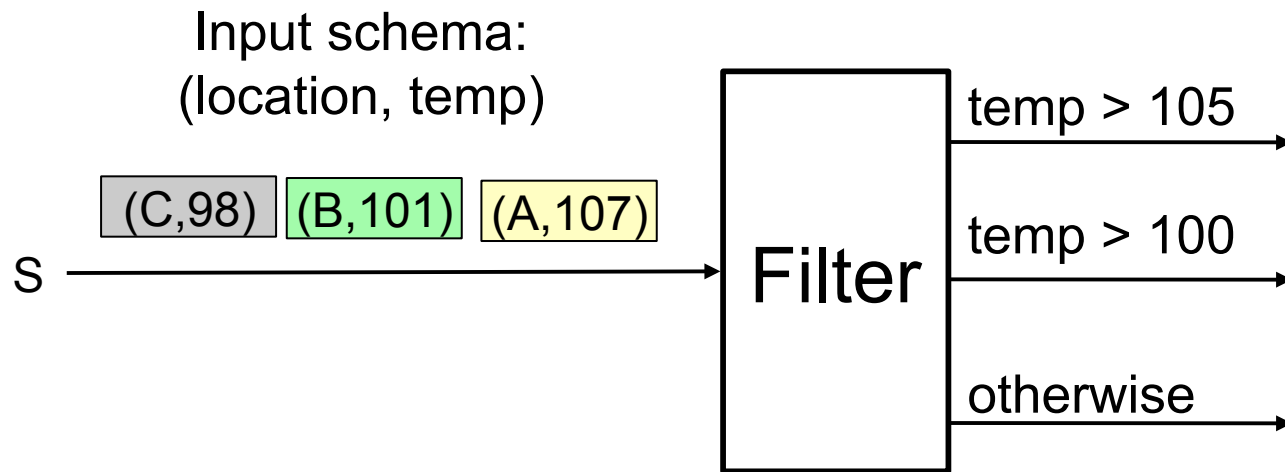- No query language

# Aurora Operators

- Order-agnostic
  - Filter
  - Map
  - Union

- Order-sensitive
  - Aggregate
  - Join
  - BSort, Resample

- **Why do we need new operators?**
  - Ops cannot block & cannot accumulate state that grows with input
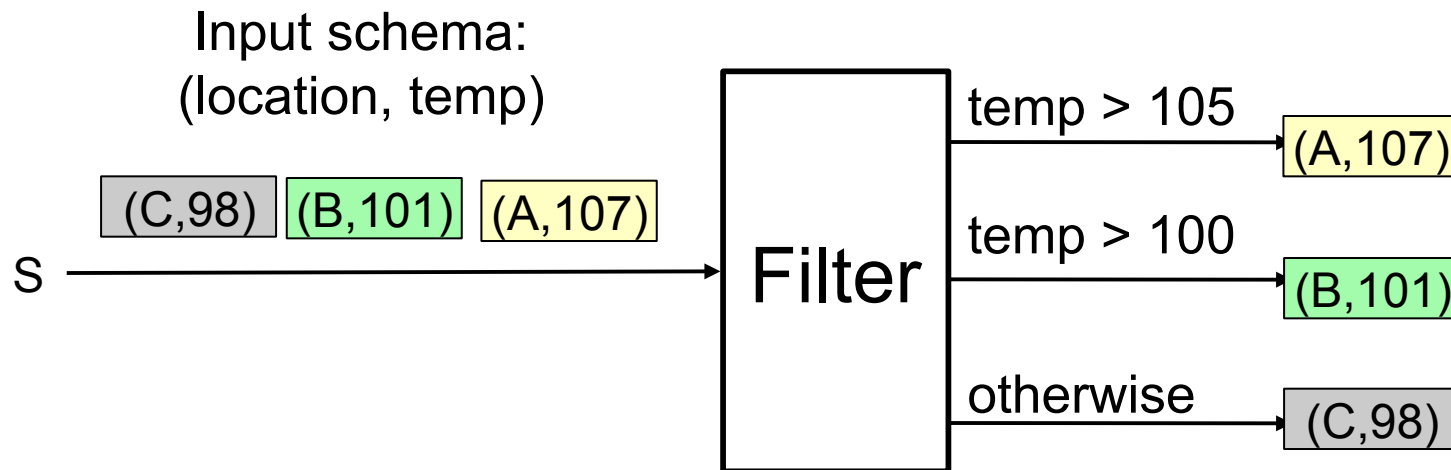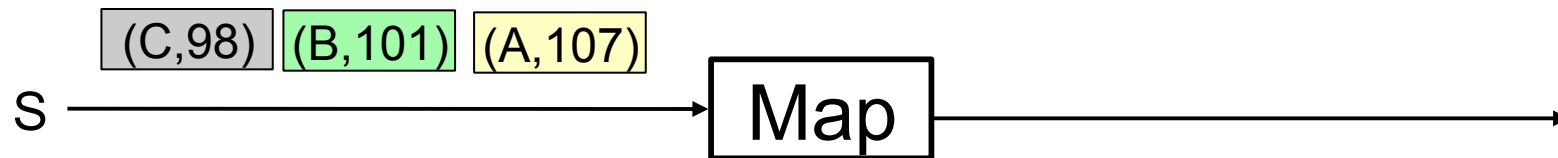
# Filter Example

Input tuples

Output tuples

Input schema:
(location, temp)

(C,98) (B,101) (A,107)

S ——→

Filter

temp > 105

temp > 100

otherwise

# Filter Example

Input tuples

Output tuples

Input schema:
(location, temp)

(C,98)  (B,101)  (A,107)

S ——→

Filter

temp > 105 → (A,107)

temp > 100 → (B,101)

otherwise → (C,98)

# Map Example

(C,98) (B,101) (A,107)

S ──────────────→ | Map | ─────────────→

new.location = old.location
new.temp_celcius = 5/9*(old.temp - 32)

# Map Example

(C,98) (B,101) (A,107)     Map     (C,37) (B,38) (A,42)

S
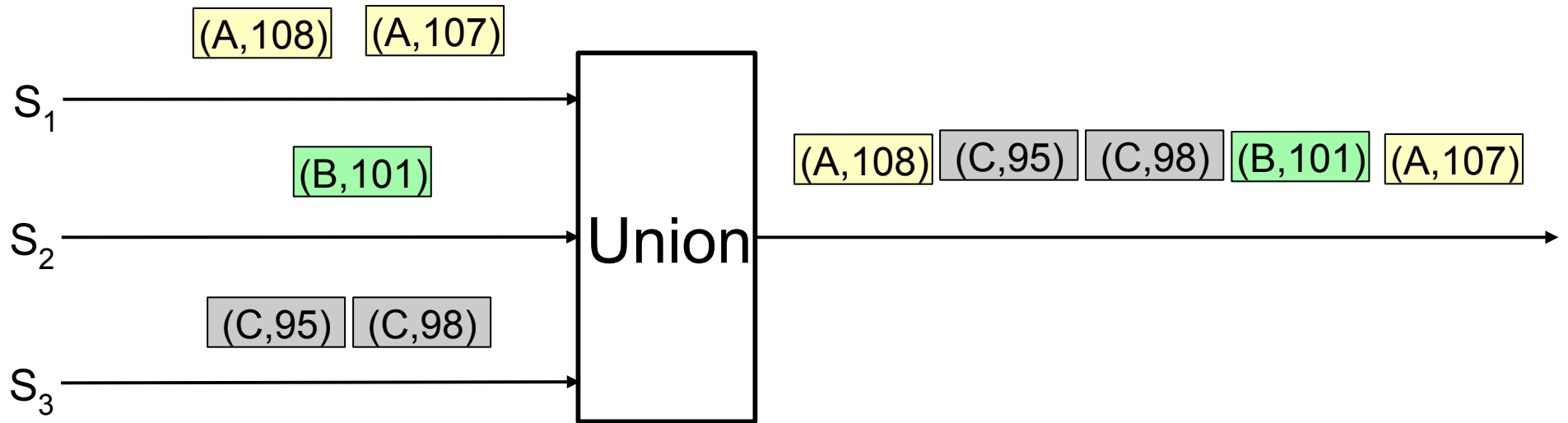
new.location = old.location
new.temp_celcius = 5/9*(old.temp - 32)

# Union Example

(A,108)   (A,107)

$S_1$ →

(B,101)

$S_2$ →   Union →

(C,95)   (C,98)

$S_3$ →

# Union Example

S₁ → (A,108) (A,107)

S₂ → (B,101)

S₃ → (C,95) (C,98)

Union → (A,108) (C,95) (C,98) (B,101) (A,107)

# Aggregate Example

S

Input tuples

(2:03,B,63)  Input schema: (time, location, temp)

(2:05,A,72)

(1:35,A,76)

(1:34,B,65)

Output tuples

(1:00,B,65)

(1:00,A,74)

Agg

Operator parameters:
- group by location
- average temp,
- order on time, window size 1 h, advance 1h

# Aggregate Example

S

Input tuples

(2:03,B,63)   Input schema: (time, location, temp)

(2:05,A,72)

(1:35,A,76)

(1:34,B,65)

(1:00,B,65)                          Output tuples

(1:00,A,74)                  (1:00,B,65)   (1:00,A,75)

Agg

Operator parameters:
- group by location
- average temp,
- order on time, window size 1 h, advance 1h

# Join Example

Input tuples

Output tuples

S

(2:05,B,72)

(1:35,A,76)

(1:00,C,74)

| (1:35,A,76, | 1:00,A,0.3) |
| (1:00,C,74, | 1:34,C,0.4) |
| (1:35,A,76, | 2:03,A,0.2) |

Join

(1:00,A,0.3)

(1:34,C,0.4)

(2:03,A,0.2)

Operator parameters:
- S order on time, R order on time, window size 1h
- predicate S.location = R.location,
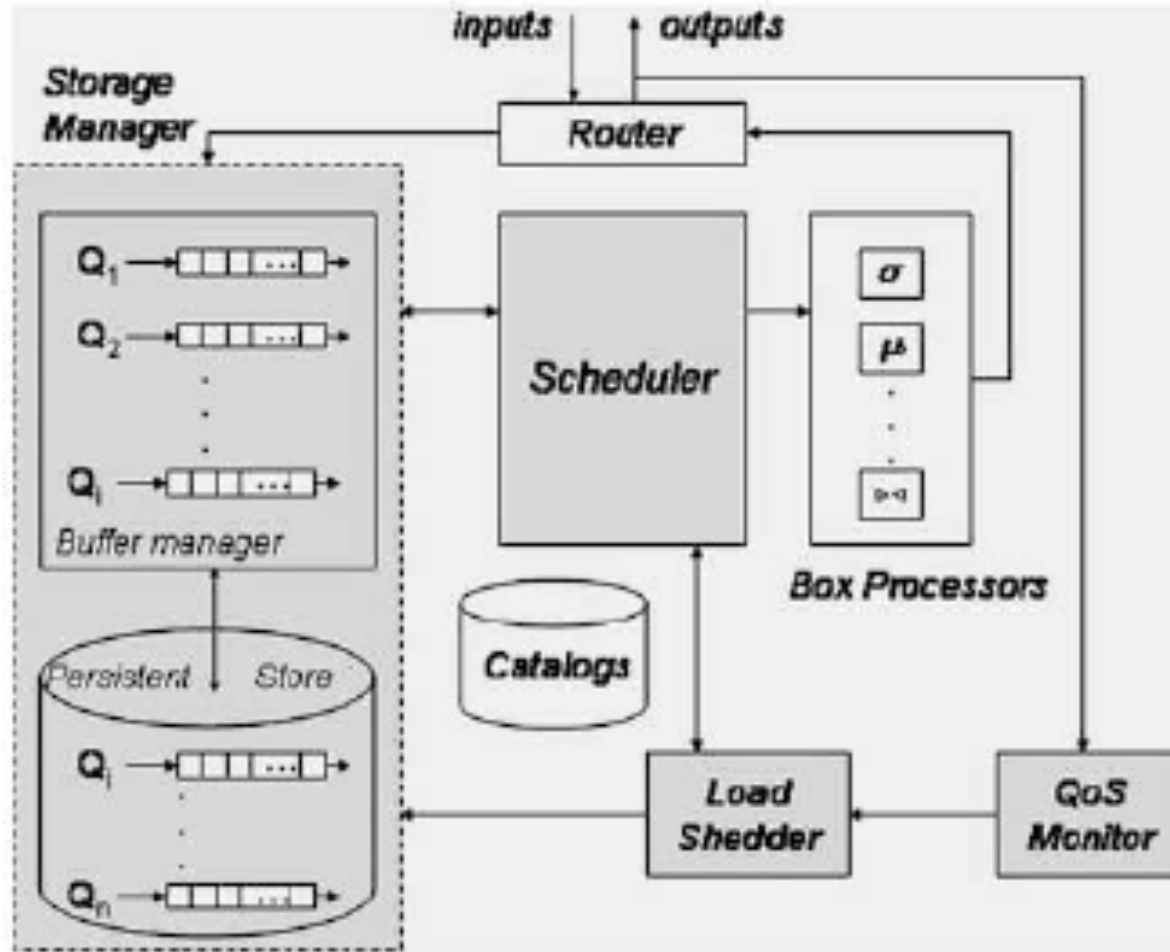
R

# Sample Query

- Application: network intrusion detection

- Schema of input stream
  `(src_ip,src_port,dst_ip,dst_port,time)`

- **Query**
  – Alert me if an IP address establishes more than 100 connections per minute
  – and within 30 seconds of that event
  – the IP tries to connect to more than 10 distinct ports within a minute

# Processing Model



[Figure 3 from Abadi:03]    CSE 544 - Winter 2009    23

# Additional Features

- **Load management**
  - What happens when system is overloaded?

- **Fault-tolerance**
  - What happens if a node fails?
  - What happens if the network fails?
  - What happens if input data is wrong?

- **Exploiting data archives**
  - Historical queries, ad-hoc queries
  - Integrating push-based processing with pull-based

- ...