

Models and Issues in Data Stream Systems

Rajeev Motwani

Stanford University

(with Brian Babcock, Shivnath Babu,
Mayur Datar, and Jennifer Widom)

STREAM Project Members: Arvind Arasu, Gurmeet Manku,
Liadan O'Callaghan, Justin Rosentstein, Qi Sun, Rohit Varma

PODS 2002

1

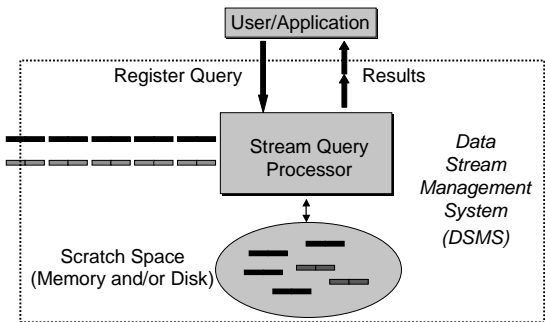
Data Streams

- Traditional DBMS – data stored in finite, persistent data sets
- New Applications – data input as continuous, ordered data streams
 - Network monitoring and traffic engineering
 - Telecom call records
 - Network security
 - Financial applications
 - Sensor networks
 - Manufacturing processes
 - Web logs and clickstreams
 - Massive data sets

PODS 2002

2

Data Stream Management System



PODS 2002

3

Meta-Questions

- Killer-apps
 - Application stream rates exceed DBMS capacity?
 - Can DSMS handle high rates anyway?
- Motivation
 - Need for general-purpose DSMS?
 - Not ad-hoc, application-specific systems?
- Non-Trivial
 - DSMS = merely DBMS with enhanced support for triggers, temporal constructs, data rate mgmt?

PODS 2002

4

Sample Applications

- Network security (e.g., iPolicy, NetForensics/Cisco, Niksun)
 - Network packet streams, user session information
 - Queries: URL filtering, detecting intrusions & DOS attacks & viruses
- Financial applications (e.g., Traderbot)
 - Streams of trading data, stock tickers, news feeds
 - Queries: arbitrage opportunities, analytics, patterns
 - SEC requirement on closing trades

PODS 2002

5

Executive Summary

- Data Stream Management Systems (DSMS)
 - Highlight issues and motivate research
 - Not a tutorial or comprehensive survey
- Caveats
 - Personal view of emerging field
 - ⊗ Stanford STREAM Project bias
 - ⊗ Cannot cover all projects in detail

PODS 2002

6

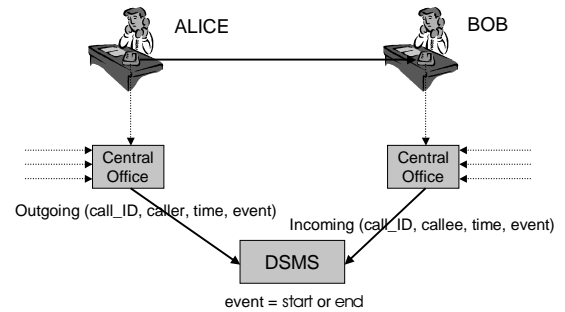
DBMS versus DSMS

- | | |
|---|---|
| • Persistent relations | • Transient streams |
| • One-time queries | • Continuous queries |
| • Random access | • Sequential access |
| • "Unbounded" disk store | • Bounded main memory |
| • Only current state matters | • History/arrival-order is critical |
| • Passive repository | • Active stores |
| • Relatively low update rate | • Possibly multi-GB arrival rate |
| • No real-time services | • Real-time requirements |
| • Assume precise data | • Data stale/imprecise |
| • Access plan determined by query processor, physical DB design | • Unpredictable/variable data arrival and characteristics |

PODS 2002

7

Making Things Concrete



PODS 2002

8

Query 1 (self-join)

- Find all outgoing calls longer than 2 minutes
- ```
SELECT O1.call_ID, O1.caller
FROM Outgoing O1, Outgoing O2
WHERE (O2.time - O1.time > 2
 AND O1.call_ID = O2.call_ID
 AND O1.event = start
 AND O2.event = end)
```
- Result requires unbounded storage
  - Can provide result as data stream
  - Can output after 2 min, without seeing end

PODS 2002

9

## Query 2 (join)

- Pair up callers and callees
- ```
SELECT O.caller, I.callee
FROM   Outgoing O, Incoming I
WHERE  O.call_ID = I.call_ID
```
- Can still provide result as data stream
 - Requires unbounded temporary storage ...
 - ... unless streams are near-synchronized

PODS 2002

10

Query 3 (group-by aggregation)

- Total connection time for each caller
- ```
SELECT O1.caller, sum(O2.time - O1.time)
FROM Outgoing O1, Outgoing O2
WHERE (O1.call_ID = O2.call_ID
 AND O1.event = start
 AND O2.event = end)
GROUP BY O1.caller
```
- Cannot provide result in (append-only) stream
    - Output updates?
    - Provide current value on demand?
    - Memory?

PODS 2002

11

## Outline of Remaining Talk

- Stream Models and DSMS Architectures
- Query Processing
- Runtime and Systems Issues
- Algorithms
- Conclusion

PODS 2002

12

## Data Model

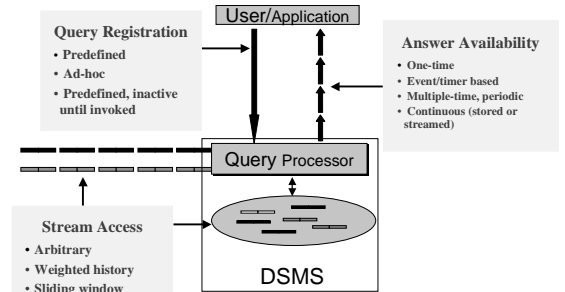
- Append-only
  - Call records
- Updates
  - Stock tickers
- Deletes
  - Transactional data
- Meta-Data
  - Control signals, punctuations

System Internals – probably need all above

PODS 2002

13

## Query Model



14

## Related Database Technology

- DSMS must use ideas, but none is substitute
  - Triggers, Materialized Views in Conventional DBMS
  - Main-Memory Databases
  - Distributed Databases
  - Pub/Sub Systems
  - Active Databases
  - Sequence/Temporal/Timeseries Databases
  - Realtime Databases
  - Adaptive, Online, Partial Results
- Novelty in DSMS
  - Semantics: input ordering, streaming output, ...
  - State: cannot store unending streams, yet need history
  - Performance: rate, variability, imprecision, ...

PODS 2002

15

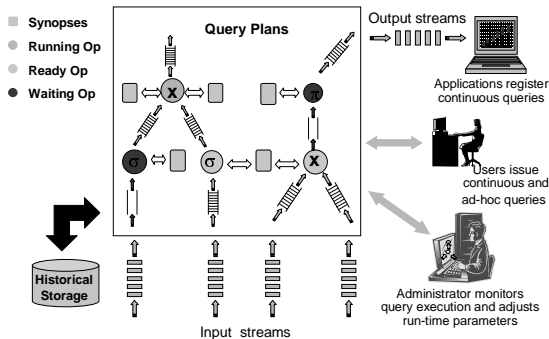
## Stream Projects

- Amazon/Cougar (Cornell) – sensors
- Aurora (Brown/MIT) – sensor monitoring, dataflow
- Hancock (AT&T) – telecom streams
- Niagara (OGI/Wisconsin) – Internet XML databases
- OpenCQ (Georgia) – triggers, incr. view maintenance
- Stream (Stanford) – general-purpose DSMS
- Tapestry (Xerox) – pub/sub content-based filtering
- Telegraph (Berkeley) – adaptive engine for sensors
- Tribeca (Bellcore) – network monitoring

PODS 2002

16

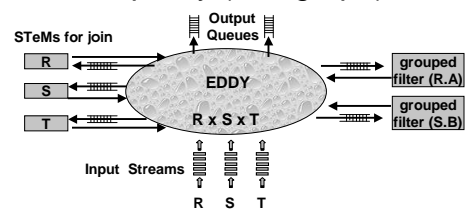
## Aurora/STREAM Overview



PODS 2002

17

## Adaptivity (Telegraph)



- Runtime Adaptivity
- Multi-query Optimization
- Framework – implements arbitrary schemes

PODS 2002

18

### Query-Split Scheme (Niagara)

Quotes.XML

trig.Act.i

scan

file i

split

Symbol = Const.Value

join

scan

file j

trig.Act.j

scan

constant table

|      |        |
|------|--------|
| ...  | ...    |
| IBM  | file i |
| MSFT | file j |
| ...  | ...    |

- Aggregate subscription for efficiency
- Split – evaluate trigger only when file updated
- Triggers – multi-query optimization

PODS 2002 19

### Shared Predicates [Niagara, Telegraph]

Predicates for R.A

- R.A > 1
- R.A > 7
- R.A > 11
- R.A < 3
- R.A < 5
- R.A = 6
- R.A = 8
- R.A = 9

Tree structure:

- Root: 7
  - Left child: 1
    - Leaf: A > 1
  - Right child: 11
    - Leaf: A > 7
    - Leaf: A > 11
- Root: 3
  - Left child: A < 3
  - Right child: A < 5
- Root: 6 | 8
- Root: 9

Output: Tuple A=8

PODS 2002 20

### Outline of Remaining Talk

- Stream Models and DSMS Architectures
- Query Processing
- Runtime and Systems Issues
- Algorithms
- Conclusion

PODS 2002 21

### Blocking Operators

- Blocking
  - No output until entire input seen
  - Streams – input never ends
- Simple Aggregates – output “update” stream
- Set Output (sort, group-by)
  - Root – could maintain output data structure
  - Intermediate nodes – try non-blocking analogs
  - Example – juggle for sort [Raman, R, Hellerstein]
  - Punctuations and constraints
- Join
  - non-blocking, but intermediate state?
  - sliding-window restrictions

PODS 2002 22

### Punctuations [Tucker, Maier, Sheard, Fegaras]

- Assertion about future stream contents
- Unblocks operators, reduces state

State/Index

R.A < 10

R.A = 10

group-by

X

R

S

P: S.A = 10

- Future Work
  - Inserted at source or internal (operator signaling)?
  - Does P unblock Q? Exists P? Rewrite Q?
  - Relation between P and memory for Q?

PODS 2002 23

### Constraints

- Schema-level: ordering, referential integrity, many-one joins
- Instance-level: punctuations
- Query-level: windowed join (nearby tuples only)

- [Babu-Widom]
  - Input – multi-stream SPJ query, schema-level constraints
  - Output – plan with low intermediate state for joins
- Future Work
  - Query-level constraints? Combining constraints?
  - Relaxed constraints (near-sorted, near-clustered)
  - Exploiting constraints in intra-operator signaling

PODS 2002 24

## Impact of Limited Memory

- Continuous streams grow unboundedly
- Queries may require unbounded memory
- [ABBMW 02]
  - a priori memory bounds for query
  - Conjunctive queries with arithmetic comparisons
  - Queries with join need domain restrictions
  - Impact of duplication elimination
- Open – general queries

PODS 2002

25

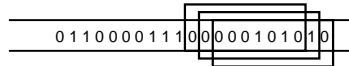
## Approximate Query Evaluation

- Why?
  - Handling load – streams coming too fast
  - Avoid unbounded storage and computation
  - Ad hoc queries need approximate history
- How? Sliding windows, synopsis, samples, load-shed
- Major Issues?
  - Metric for set-valued queries
  - Composition of approximate operators
  - How is it understood/controlled by user?
  - Integrate into query language
  - Query planning and interaction with resource allocation
  - Accuracy-efficiency-storage tradeoff and global metric

PODS 2002

26

## Sliding Window Approximation



- Why?
  - Approximation technique for bounded memory
  - Natural in applications (emphasizes recent data)
  - Well-specified and deterministic semantics
- Issues
  - Extend relational algebra, SQL, query optimization
  - Algorithmic work
  - Timestamps?

PODS 2002

27

## Timestamps

- Explicit
  - Injected by data source
  - Models real-world event represented by tuple
  - Tuples may be out-of-order, but if near-ordered can reorder with small buffers
- Implicit
  - Introduced as special field by DSMS
  - Arrival time in system
  - Enables order-based querying and sliding windows
- Issues
  - Distributed streams?
  - Composite tuples created by DSMS?

PODS 2002

28

## Timestamps in JOIN Output



### Approach 1

- User-specified, with defaults
- Compute output timestamp
- Must output in order of timestamps
- Better for Explicit Timestamp
- Need more buffering
- Get precise semantics and user-understanding

### Approach 2

- Best-effort, no guarantee
- Output timestamp is exit-time
- Tuples arriving earlier more likely to exit earlier
- Better for Implicit Timestamp
- Maximum flexibility to system
- Difficult to impose precise semantics

PODS 2002

29

## Approximate via Load-Shedding

Handles scan and processing rate mismatch

### Input Load-Shedding

- Sample incoming tuples
- Use when scan rate is bottleneck
- Positive – online aggregation [Hellerstein, Haas, Wang]
- Negative – join sampling [Chaudhuri, Motwani, Narasaya]

### Output Load-Shedding

- Buffer input infrequent output
- Use when query processing is bottleneck
- Example – XJoin [Urhan, Franklin]
- Exploit synopses

PODS 2002

30

## Distributed Query Evaluation

- Logical stream = many physical streams
  - maintain top 100 Yahoo pages
- Correlate streams at distributed servers
  - network monitoring
- Many streams controlled by few servers
  - sensor networks
- Issues
  - Move processing to streams, not streams to processors
  - Approximation-bandwidth tradeoff

PODS 2002

31

## Example: Distributed Streams

- Maintain top 100 Yahoo pages
  - Pages served by geographically distributed servers
  - Must aggregate server logs
  - Minimize communication
- Pushing processing to streams
  - Most pages not in top 100
  - Avoid communicating about such pages
  - Send updates about relevant pages only
  - Requires server coordination

PODS 2002

32

## Stream Query Language?

- SQL extension
- Sliding windows as first-class construct
  - Awkward in SQL, needs reference to timestamps
  - SQL-99 allows aggregations over sliding windows
- Sampling/approximation/load-shedding/QoS support?
- Stream relational algebra and rewrite rules
  - Aurora and STREAM
  - Sequence/Temporal Databases

PODS 2002

33

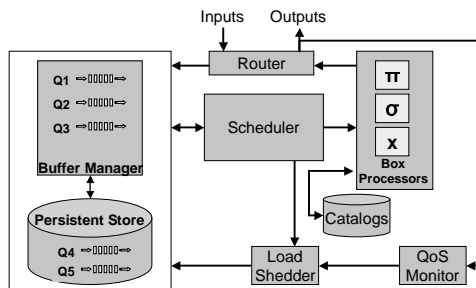
## Outline of Remaining Talk

- Stream Models and DSMS Architectures
- Query Processing
- Runtime and Systems Issues
- Algorithms
- Conclusion

PODS 2002

34

## Aurora Run-time Architecture



PODS 2002

35

## DSMS Internals

- Query plans: operators, synopses, queues
- Memory management
  - Dynamic Allocation – queries, operators, queues, synopses
  - Graceful adaptation to reallocation
  - Impact on throughput and precision
- Operator scheduling
  - Variable-rate streams, varying operator/query requirements
  - Response time and QoS
  - Load-shedding
  - Interaction with queue/memory management

PODS 2002

36

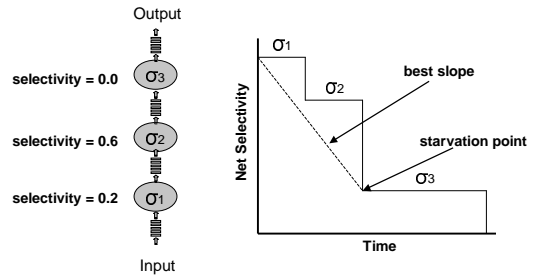
## Queue Memory and Scheduling [Babcock, Babu, Datar, Motwani]

- Goal
  - Given – query plan and selectivity estimates
  - Schedule – tuples through operator chains
- Minimize total queue memory
  - Best-slope scheduling is near-optimal
  - Danger of starvation for some tuples
- Minimize tuple response time
  - Schedule tuple completely through operator chain
  - Danger of exceeding memory bound
- Open – graceful combination and adaptivity

PODS 2002

37

## Queue Memory and Scheduling [Babcock, Babu, Datar, Motwani]



PODS 2002

38

## Precision-Resource Tradeoff

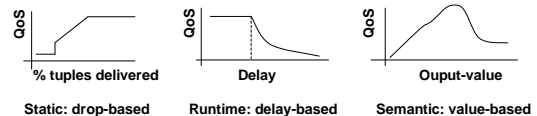
- Resources – memory, computation, I/O
- Global Optimization Problem
  - Input: queries with alternate plans, importance weights
  - Precision: function of resource allocation to queries/operators
  - Goal: select plans, allocate resources, maximize precision
- Memory Allocation Algorithm [Varma, Widom]
  - Model – single query plan, simple precision model
  - Rules for precision of composed operators
  - Non-linear numerical optimization formulation
- Open – Combinatorial algorithm? General case?

PODS 2002

39

## Rate-Based & QoS Optimization

- [Viglas, Naughton]
  - Optimizer goal is to increase throughput
  - Model for output-rates as function of input-rates
  - Designing optimizers?
- Aurora – QoS approach to load-shedding



PODS 2002

40

## Outline of Remaining Talk

- Stream Models and DSMS Architectures
- Query Processing
- Runtime and Systems Issues
- Algorithms
- Conclusion

PODS 2002

41

## Synopses

- Queries may access or aggregate past data
- Need bounded-memory history-approximation
- Synopsis?
  - Succinct summary of old stream tuples
  - Like indexes/materialized-views, but base data is unavailable
- Examples
  - Sliding Windows
  - Samples
  - Sketches
  - Histograms
  - Wavelet representation

PODS 2002

42

### Model of Computation

Increasing time

Data Stream

Synopses/Data Structures

Memory:  $\text{poly}(1/\epsilon, \log N)$   
 Query/Update Time:  $\text{poly}(1/\epsilon, \log N)$   
 N: # tuples so far, or window size  
 $\epsilon$ : error parameter

PODS 2002 43

### Sketching Techniques

- [Alon, Matias, Szegedy] frequency moments
- [Feigenbaum et al, Indyk] extended to  $L_p$  norm
- [Dobra et al] complex aggregates over joins
- Key Subproblem – Self-Join Size Estimation
  - Stream of values from  $D = \{1, 2, \dots, N\}$
  - Let  $f_i$  = frequency of value  $i$
  - Self-join size  $S = \sum f_i^2$
  - Question – estimating  $S$  in small space?

PODS 2002 44

### Self-Join Size Estimation

- AMS Technique (randomized sketches)
  - Given  $(f_1, f_2, \dots, f_N)$
  - $Z_i = \text{random}\{-1, 1\}$
  - $X = \sum f_i Z_i$  ( $X$  incrementally computable)
- Theorem  $\text{Exp}[X^2] = \sum f_i^2$ 
  - Cross-terms  $f_i Z_i f_j Z_j$  have 0 expectation
  - Square-terms  $f_i Z_i f_i Z_i = f_i^2$
- Space =  $\log(N + \sum f_i)$
- Independent samples  $X_k$  reduce variance

PODS 2002 45

### Sliding Window Computations

[Datar, Gionis, Indyk, Motwani]

- Goal: statistics/queries
- Memory:  $o(N)$ , preferably  $\text{poly}(1/\epsilon, \log N)$
- Problem: count/sum/variance, histogram, clustering, ...
- Sample Results:  $(1+\epsilon)$ -approximation
  - Counting: Space  $O(1/\epsilon \log N)$  bits, Time  $O(1)$  amortized
  - Sum over  $[0, R]$ : Space  $O(1/\epsilon \log N (\log N + \log R))$  bits, Time  $O(\log R / \log N)$  amortized
  - $L_p$  sketches: maintain with  $\text{poly}(1/\epsilon, \log N)$  space overhead
  - Matching space lower bounds

PODS 2002 46

### Sliding Window Histograms

- Key Subproblem – Counting 1's in bit-stream
- Goal – Space  $O(\log N)$  for window size  $N$
- Problem – Accounting for expiring bits
- Idea
  - Partition/track buckets of known count
  - Error in oldest bucket only
  - Future 0's?

PODS 2002 47

### Exponential Histograms

- Buckets of exponentially increasing size
- Between  $K/2$  and  $K/2+1$  buckets of each size
  - $K = 1/\epsilon$  and  $\epsilon$  = relative error

PODS 2002 48

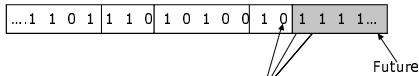


## Exponential Histograms

- Buckets of exponentially increasing size
- Between  $K/2$  and  $K/2+1$  buckets of each size
  - $K = 1/\epsilon$  and  $\epsilon =$  relative error

$K = 2$

Bucket sizes = 4, 2, 2, 1, 1, 1.



Element arrived this step.

$$C_{i-1} + C_{i-2} + \dots + C_2 + C_1 + 1 \geq (K/2) C_i$$

PODS 2002

49

## Many other results ...

- Histograms
  - V-Opt Histograms  
[Gilbert, Guha, Indyk, Kotidis, Muthukrishnan, Strauss], [Indyk]
  - End-Biased Histograms (Iceberg Queries)  
[Manku, Motwani], [Fang, Shiva, Garcia-Molina, Motwani, Ullman]
  - Equi-Width Histograms (Quantiles)  
[Manku, Rajagopalan, Lindsay], [Khanna, Greenwald]
  - Wavelets  
Seminal work [Vitter, Wang, Iyer] + many others!
- Data Mining
  - Stream Clustering  
[Guha, Mishra, Motwani, O'Callaghan]  
[O'Callaghan, Meyerson, Mishra, Guha, Motwani]
  - Decision Trees  
[Domingos, Hulten], [Domingos, Hulten, Spencer]

PODS 2002

50

## Conclusion: Future Work

- Query Processing
  - Stream Algebra and Query Languages
  - Approximations
  - Blocking, Constraints, Punctuations
- Runtime Management
  - Scheduling, Memory Management, Rate Management
  - Query Optimization (Adaptive, Multi-Query, Ad-hoc)
  - Distributed processing
- Synopses and Algorithmic Problems
- Systems
  - UI, statistics, crash recovery and transaction management
  - System development and deployment

PODS 2002

51

Thank You!

PODS 2002

52