# Adaptive Aggregation on Chip Multiprocessors

## John Cieslewicz - Kenneth A. Ross

### Columbia University

# Motivation

- Aggregation is a well understood database operation

- Seemingly easy to implement

- Chip multiprocessors introduce new challenges and opportunities

- Picking the wrong aggregation technique can result in a <u>performance penalty of more than an order of magnitude</u>
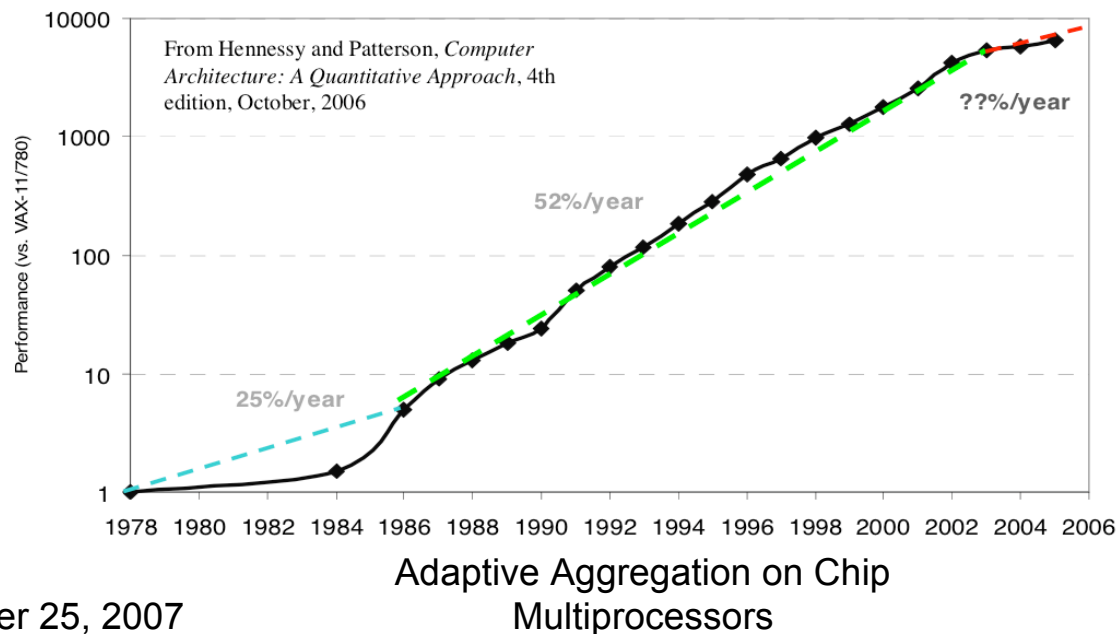
# Outline

- **Chip Multiprocessors**

- **Aggregation Strategies**

- **Modeling Performance**

- **Adaptive Aggregation**

- **Experimental Results**

Adaptive Aggregation on Chip
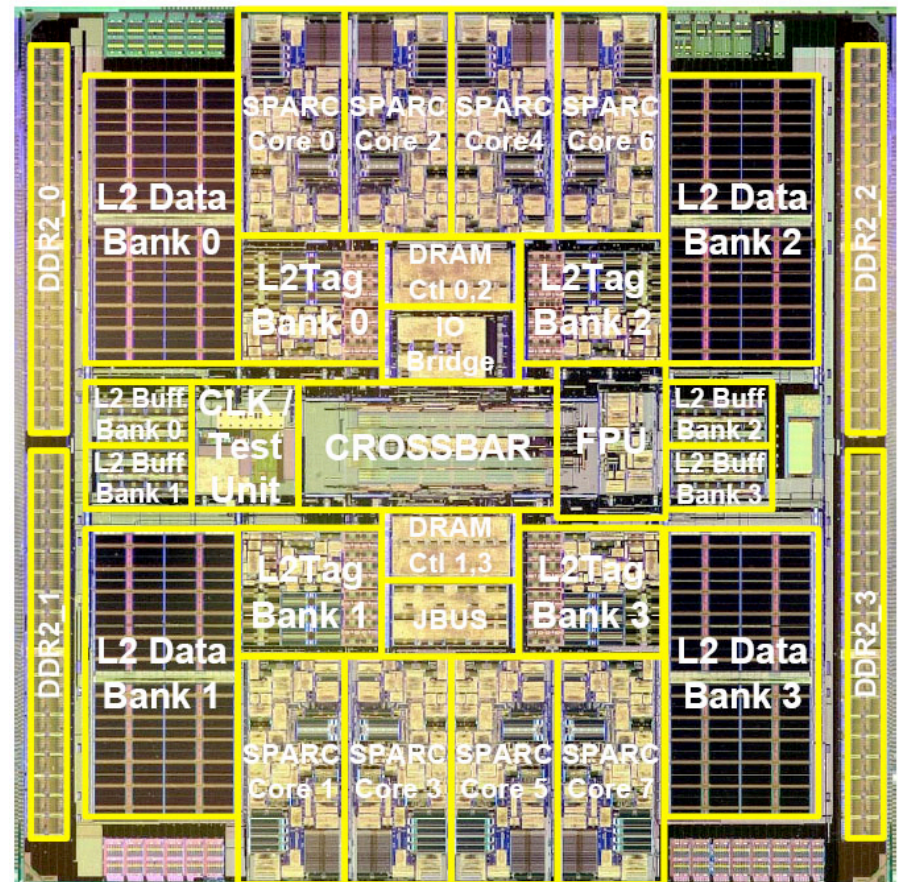Multiprocessors

# The Multicore Future

- ILP is tapped out
- Heat dissipation and power consumption problems

- Thread Level Parallelism (TLP) is the future of performance gains



Performance (vs. VAX-11/780)

From Hennessy and Patterson, *Computer Architecture: A Quantitative Approach*, 4th edition, October, 2006

25%/year

52%/year

??%/year

Adaptive Aggregation on Chip Multiprocessors

# Sun UltraSPARC T1

- 1 GHz
- 8 cores
- 4 threads / core
- 8 KB L1 D$ / core
- 16 KB L1 I$ / core
- 3 MB Shared L2
- Simple cores

Adaptive Aggregation on Chip Multiprocessors
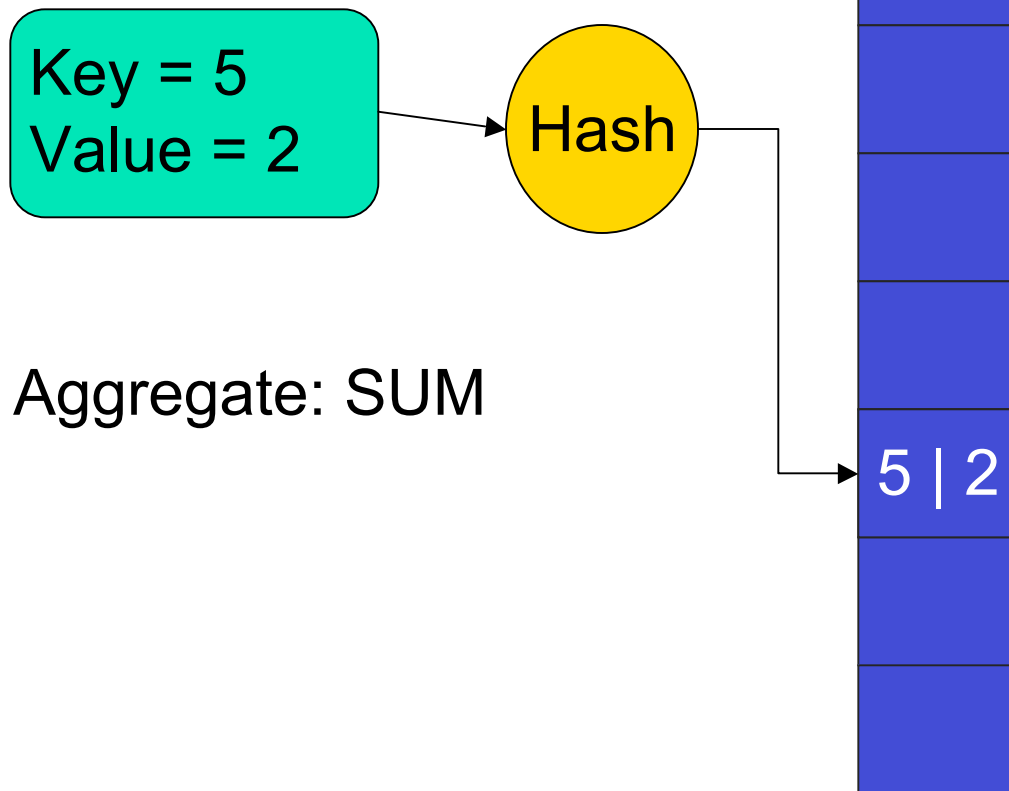
# Aggregation Overview

- Group tuples by zero or more attributes
- Compute an aggregate for each group
  - SQL standards: COUNT, SUM, AVERAGE, MIN, MAX
- Two common strategies:
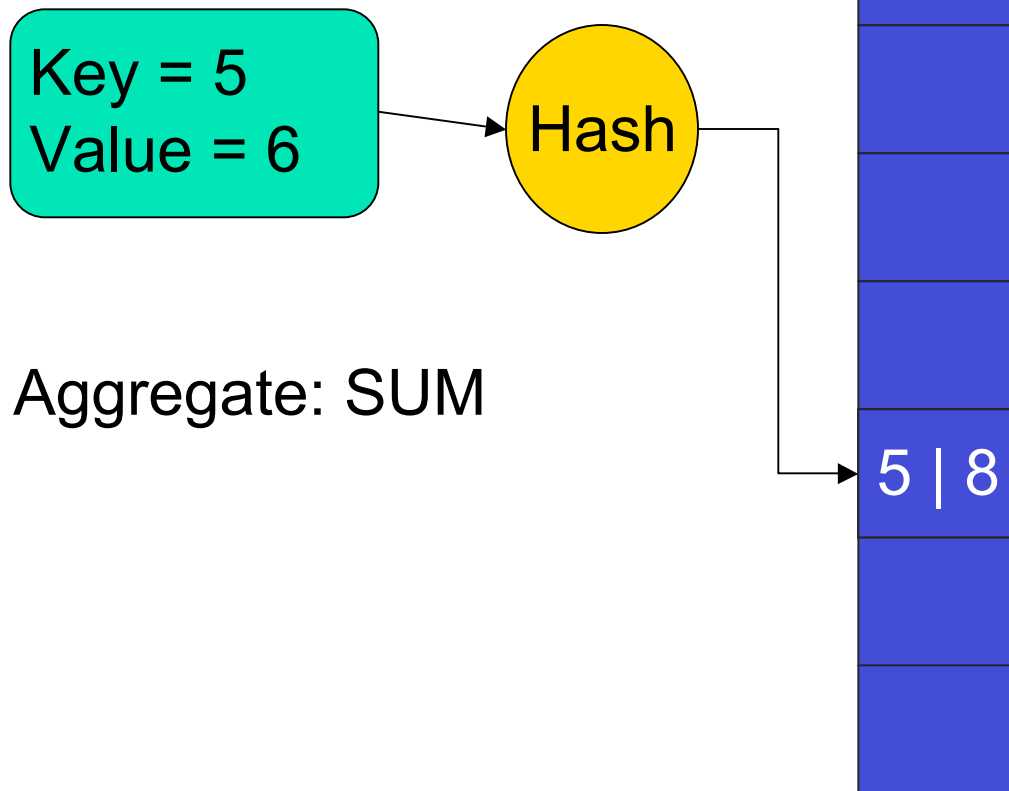  - Sorting
  - Hashing

# Hashing, not sorting

- Sort based aggregation is a blocking operation
  - All input must be materialized
- Hashing allows for better pipelining and arbitrary partitioning of the input
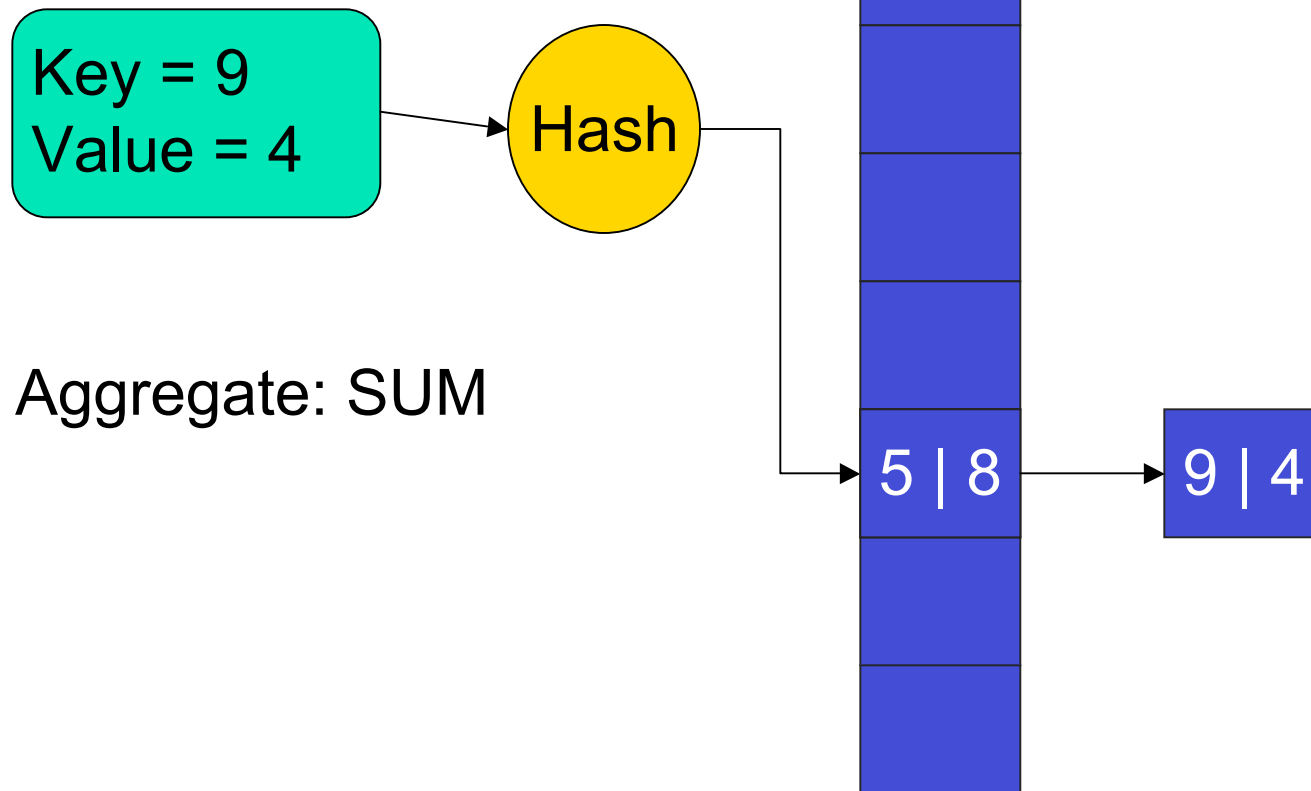- We focus only on hashing for the rest of the presentation

Adaptive Aggregation on Chip Multiprocessors

# Hash Aggregation

Key = 5
Value = 2

Hash

5 | 2

Aggregate: SUM

Adaptive Aggregation on Chip Multiprocessors

# Hash Aggregation

Key = 5
Value = 6

Hash

Aggregate: SUM

5 | 8

Adaptive Aggregation on Chip Multiprocessors

# Hash Aggregation

Key = 9
Value = 4

Hash

Aggregate: SUM

5 | 8    9 | 4

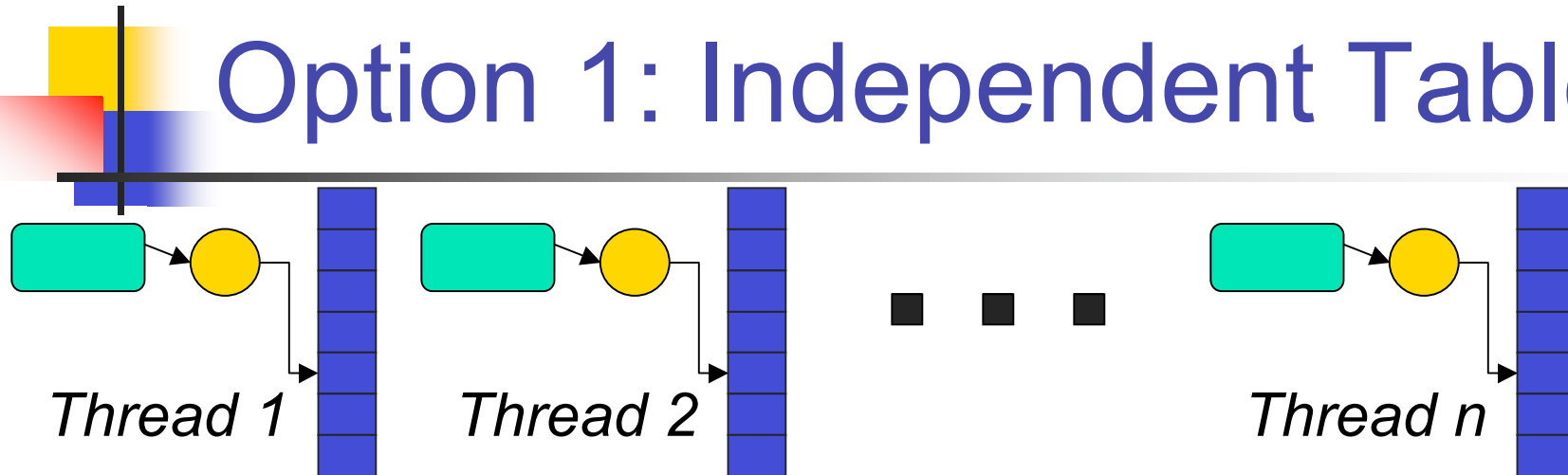Adaptive Aggregation on Chip
Multiprocessors
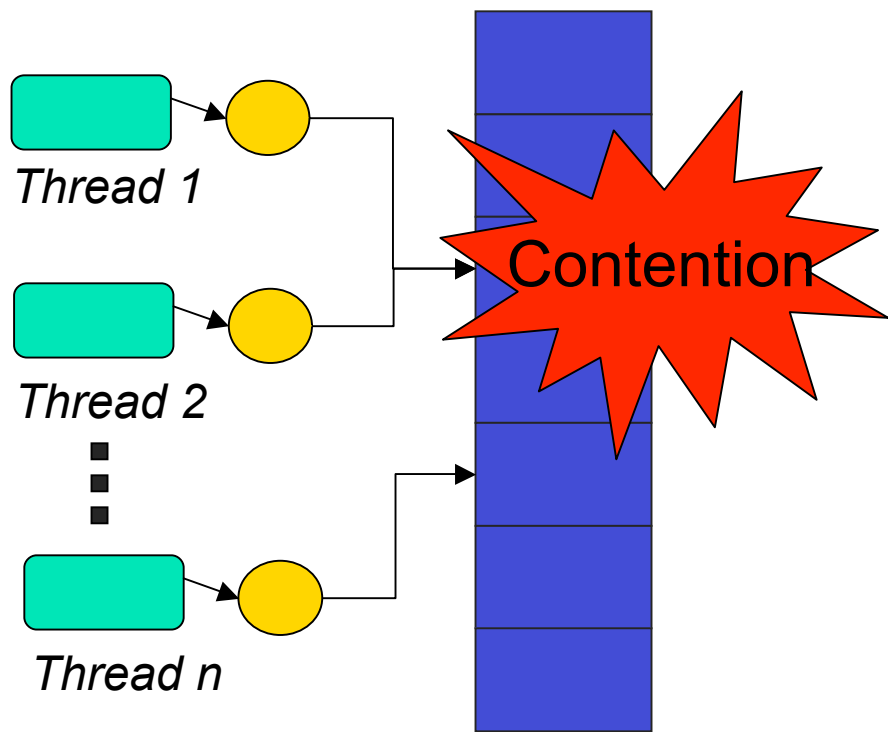
# Aggregation Implementation

- All threads share input stream
  - Read contiguous chunks
- Execute same operation
- Intra-operator sharing and conflicts are easier to reason about than inter-operator
- Instructions shared by threads
  - Instruction cache misses are expensive

# Option 1: Independent Tables



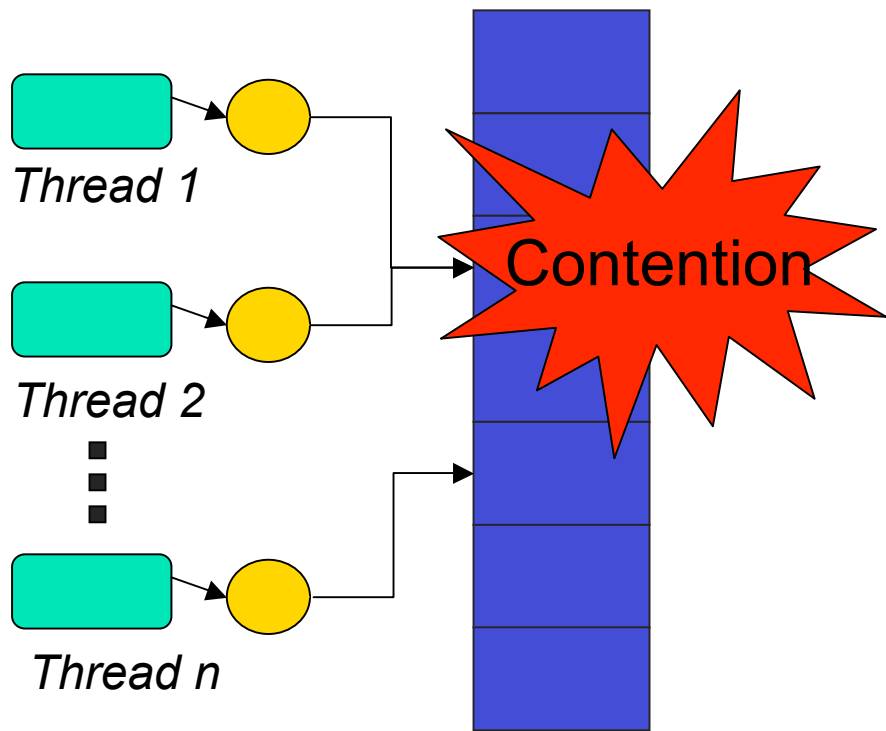*Thread 1*      *Thread 2*      *Thread n*

- Each thread has its own hash table
- Advantages:
  - No coordination between threads
- Disadvantages:
  - No sharing
  - *Capacity* and *conflict* cache misses
  - Huge memory requirement

Adaptive Aggregation on Chip
Multiprocessors

# Option 2: Global Tables with Mutexes

**Thread 1**
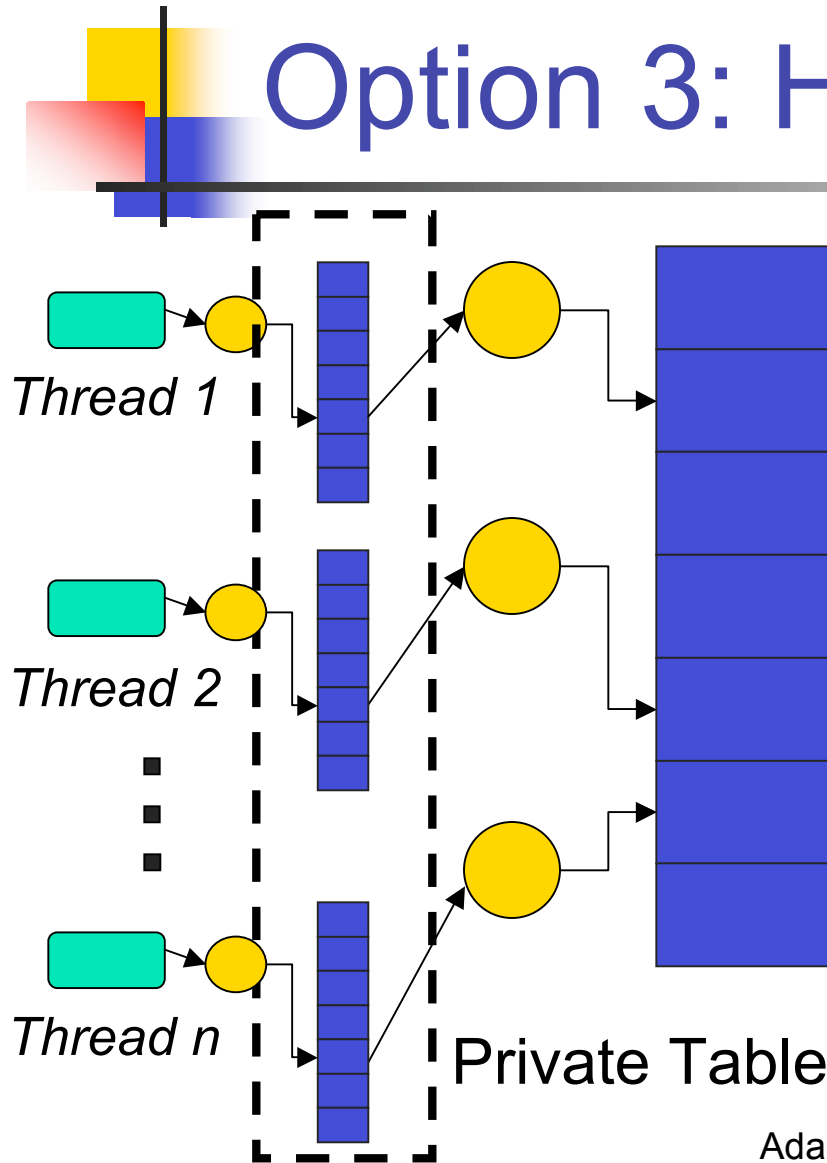
**Thread 2**

**Thread n**

Contention

- Threads share 1 table
- Advantage:
  - Shared table means more unique aggregate values fit in the cache
- Disadvantage:
  - Hash buckets must be locked to prevent race conditions
  - Contention for common keys

# Option 2: Global Tables with Atomic Instructions



- Atomicity like a mutex, but...
- No locking, use *atomic operations* for updates
- Provided by many microarchitectures
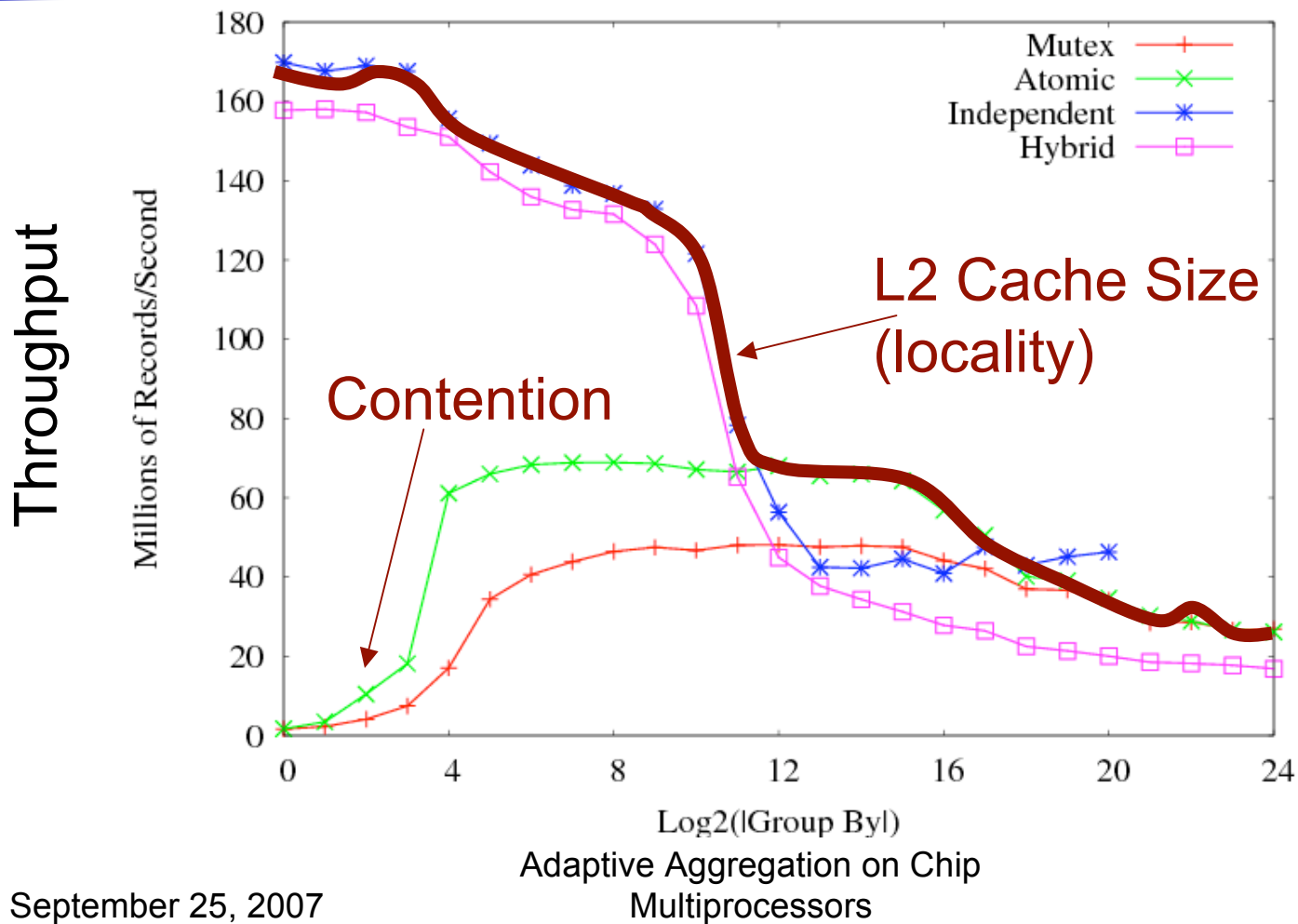- More efficient than locking, longer latency than comparable non-atomic operation

Adaptive Aggregation on Chip Multiprocessors

# Option 3: Hybrid

Independent, fixed size tables fit in L2 cache

"Spill" to global table

Advantage:
- locality
- contention free
- lower memory needs

Disadvantage
- Pure overhead if global table is better

*Thread 1*

*Thread 2*

*Thread n*

Private Tables Fit in L2

Adaptive Aggregation on Chip Multiprocessors

# Aggregation Performance

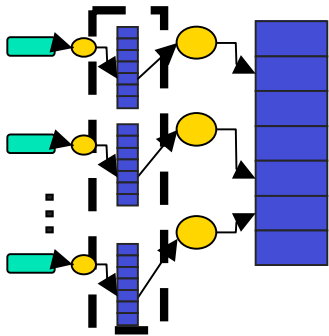Adaptive Aggregation on Chip
Multiprocessors

# What influences performance?

- ## Runs
  - Consecutive tuples with same key can be aggregated directly
- ## Locality
  - If keys repeat with temporal locality, bucket will be in the cache
- ## Contention
  - If keys repeat too often in multiple threads, contention may occur for shared hash buckets

# Modeling Performance

- Sample the input stream

- Add statistics gathering to hybrid approach

- Each thread gathers independent statistics

  - No coordination overhead

  - Each thread proceeds as fast as possible

  - Local decision may not be globally optimal
    E.g., "If every thread saw input like mine, there would be contention."

Adaptive Aggregation on Chip
Multiprocessors

# Why can't the optimizer choose?

- Statistics might be wrong or inadequate
- Static choice cannot adapt to change in distribution
- Multiple operators to choose from versus one that works well; reduces plan space
- Aggregation occurs late in plans, other operators may have introduced skew

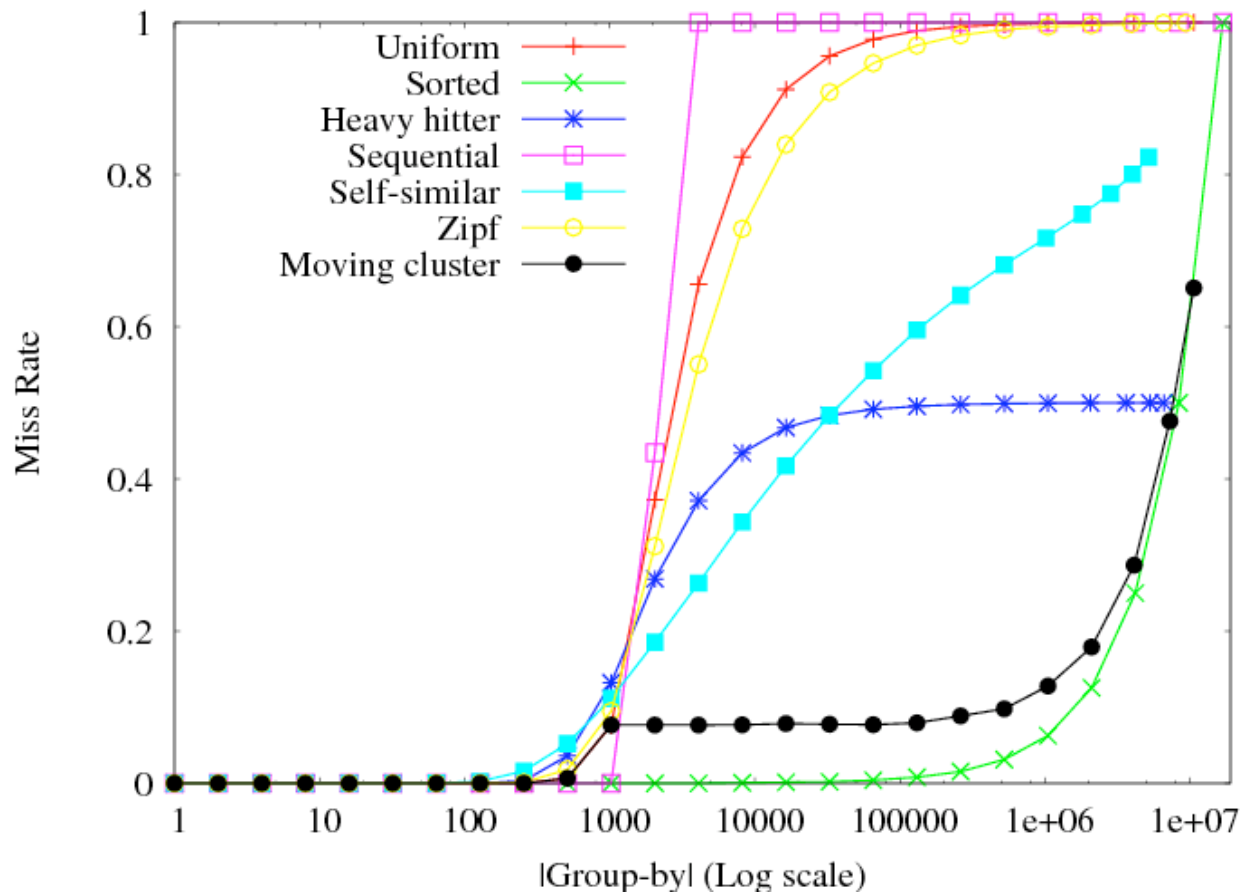Adaptive Aggregation on Chip Multiprocessors

# Sampling for Runs

- Count the number of runs seen during the sampling window, find average run length

- On uniform input, run optimization is beneficial up to |Group By| = 8

- Expected run length=$1+(1/8)^2+(1/8)^3...=8/7$

- Use run optimization if average run length exceeds 8/7

Adaptive Aggregation on Chip Multiprocessors

# Sampling for Locality

- Count "hits" in the local table
  - A "hit" is when a key is found (no insertion)
  - Tables sized to fit in L2 (likely to be a $ hit)
- Avoid compulsory misses with "warm-up"
- Locality if miss rate is less than 50%
  - Derived by the relative cost of processing a tuple in the local table compared to the cost of using a global table

Adaptive Aggregation on Chip Multiprocessors

# Miss rate

Adaptive Aggregation on Chip
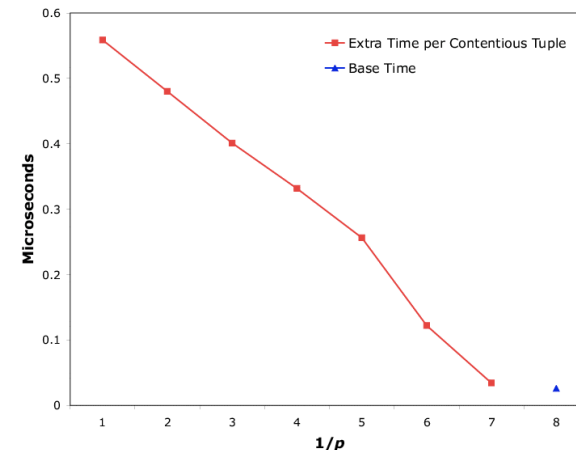Multiprocessors

# Contention

- Often subsumed by locality, except...

- Distributions with "heavy hitters" have contention without locality

- Global table must be avoided if there is contention because overhead dominates execution time

Adaptive Aggregation on Chip Multiprocessors

# Modeling Contention

- Contention directly related to the proportion of the input with the same key.
  - Contention negligible below a measurable threshold
- See paper for model of contention
  - When contention is present, the penalty per contentious tuple is *linearly related* to the inverse of the key's frequency in the input

Adaptive Aggregation on Chip Multiprocessors

# Sampling for Contention

- Count accesses to each hash bucket
- If any bucket's access count exceeds a threshold, mark it as potentially contentious[*]
- Calculate the penalty due to contention of all marked contentious buckets[*]
- If the cumulative contention penalty is sufficiently high, flag the input as contentious[*]
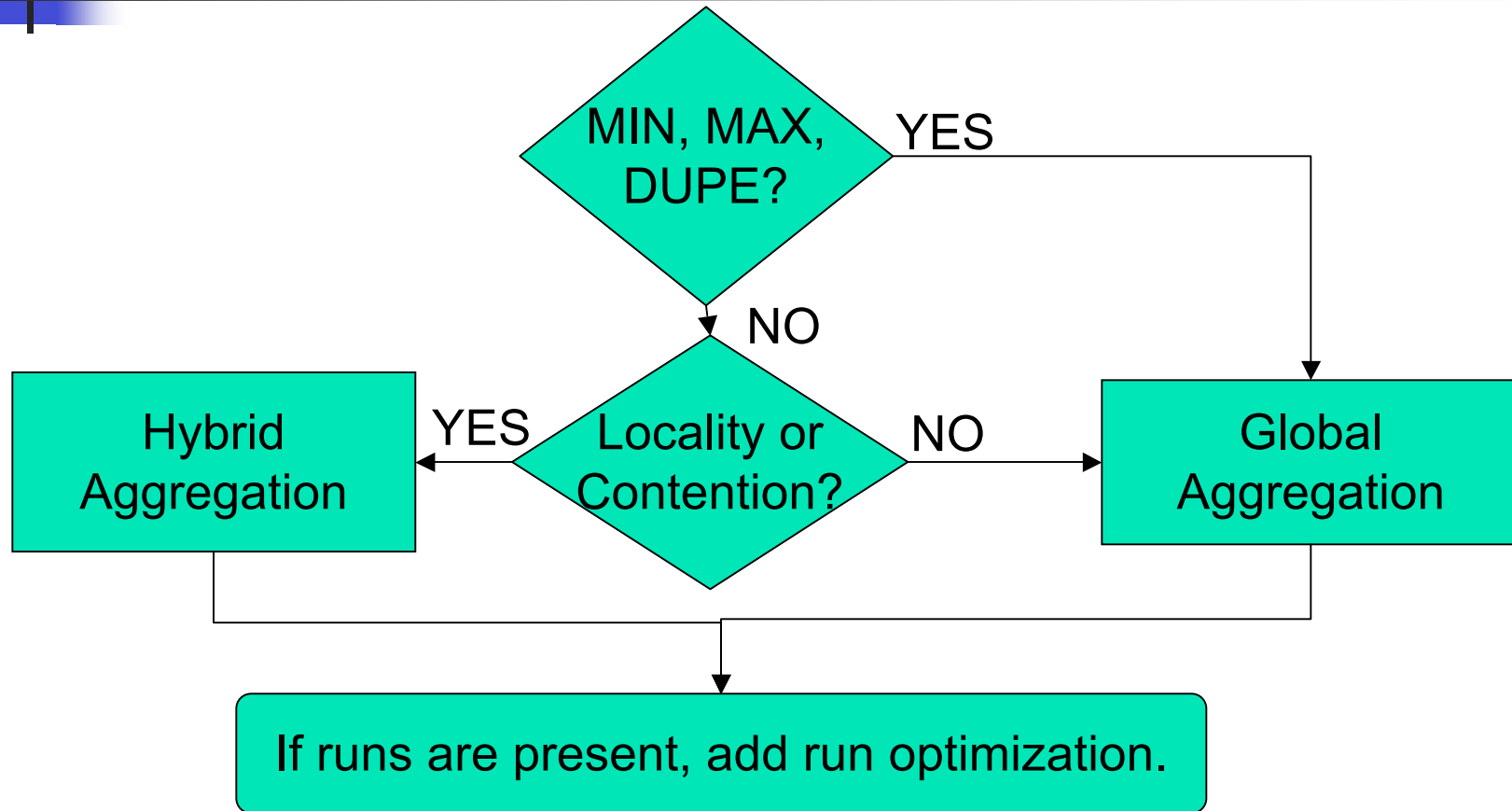
[*]See paper for a full description

# MIN, MAX, Duplicate Elimination

- **Contention Free**
- **Why? Answer: Updates are rare**
  - E.g., given a uniform input, after 99 inputs the running minimum is in the first percentile. The chance that the 100th value will update the minimum is 1%
- **Adversarial distributions exist, but can be handled with randomization**
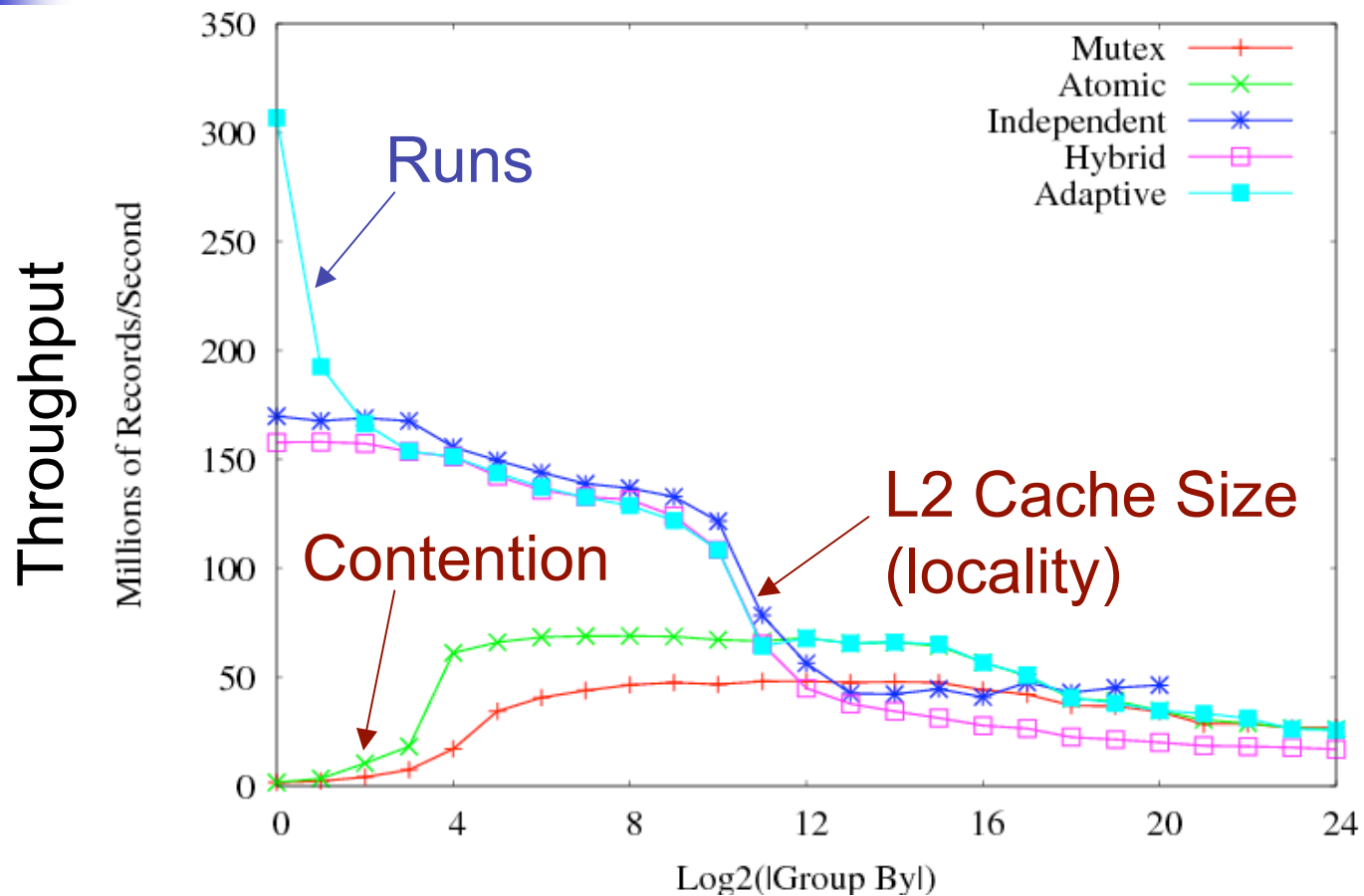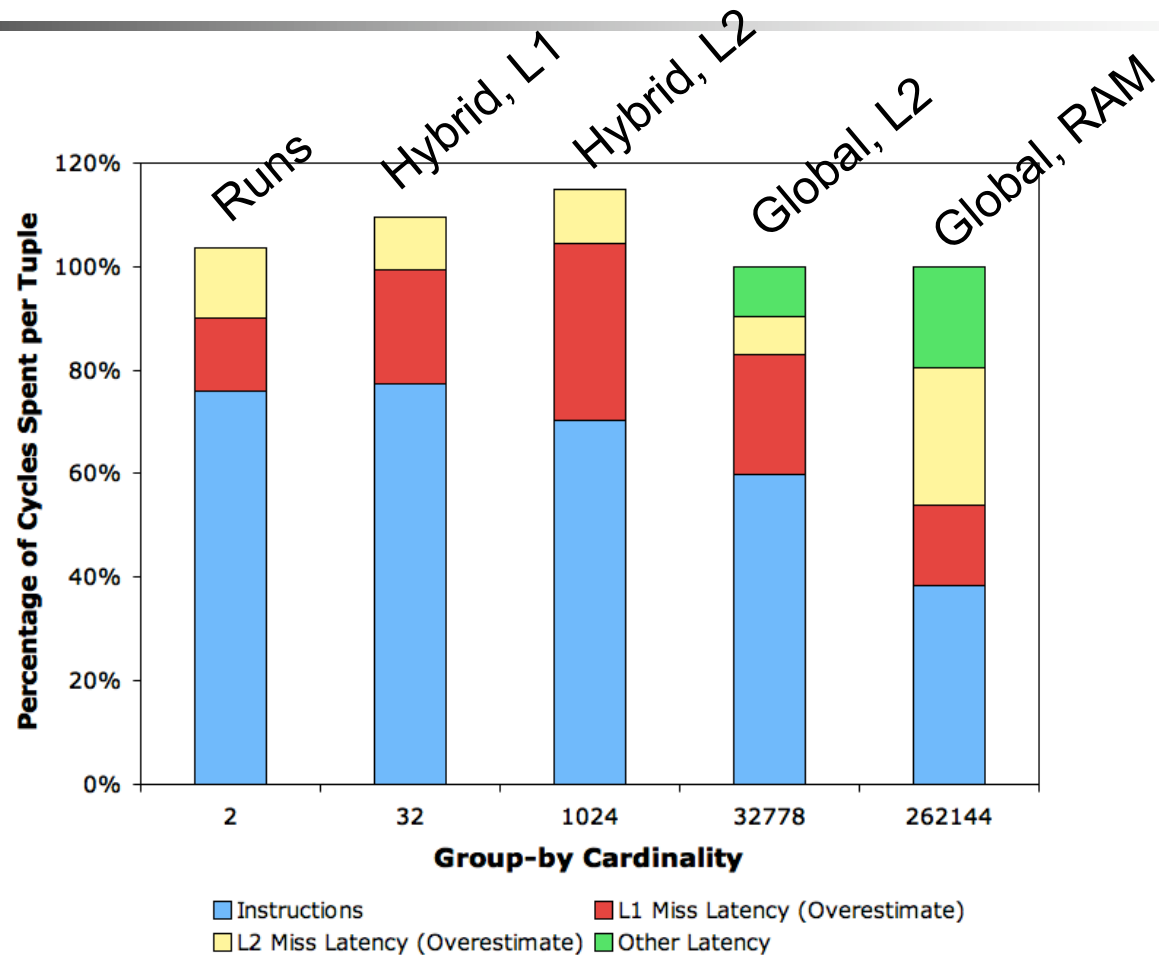
Adaptive Aggregation on Chip Multiprocessors

# Adaptive Aggregation



September 25, 2007

# Experiments

- $2^{24}$ ≈ 16 Million Input Tuples

- 7 Input Distributions, 3 Queries

- **Q1:** `SELECT G, count(*), sum(V), sum(V*V)`
  `FROM R GROUP BY G`

- **Q2:** `SELECT G, max(V), min(V), max(V)`
  `FROM R GROUP BY G`
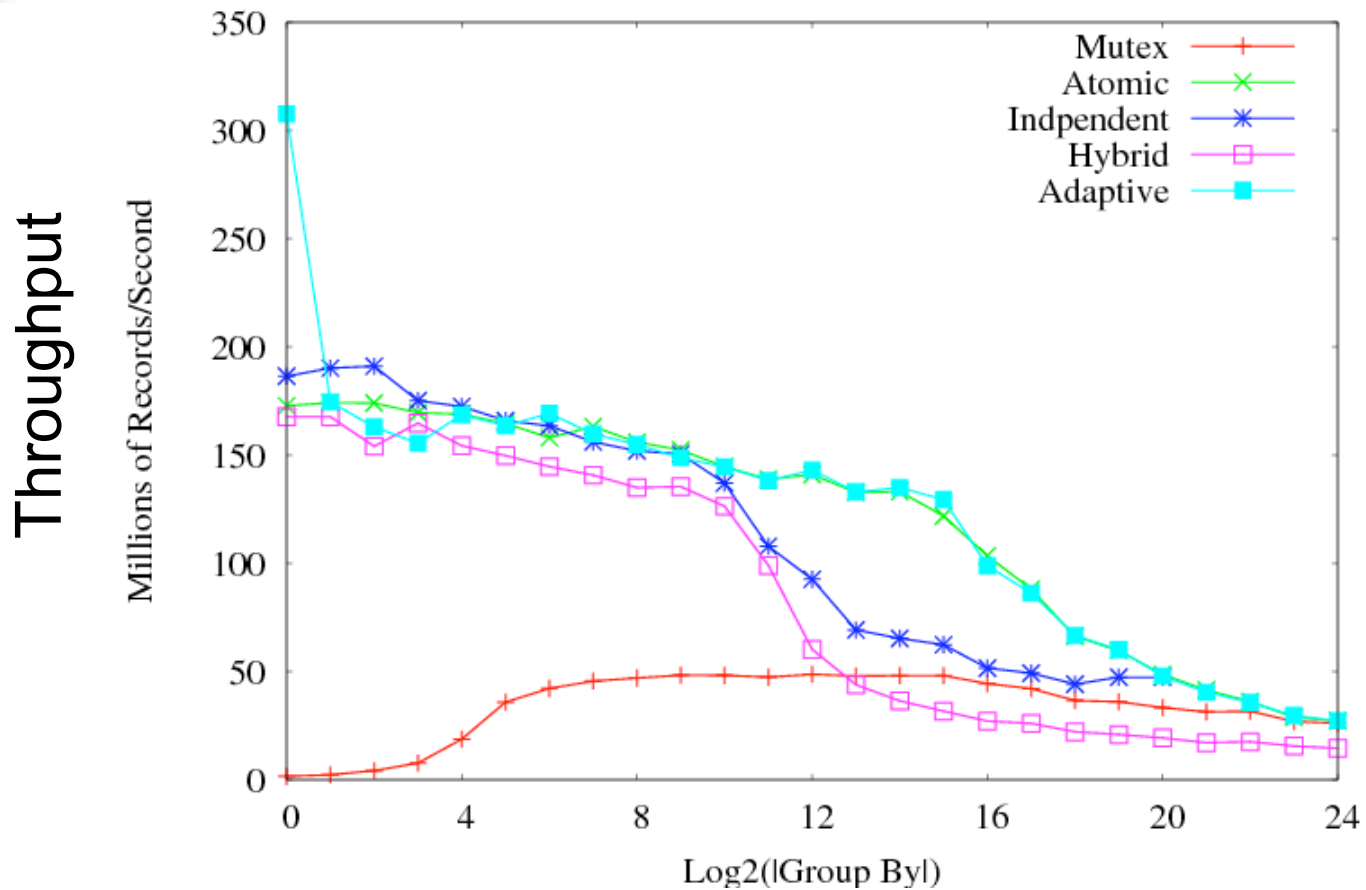
- **Q3:** `SELECT DISTINCT G FROM R`
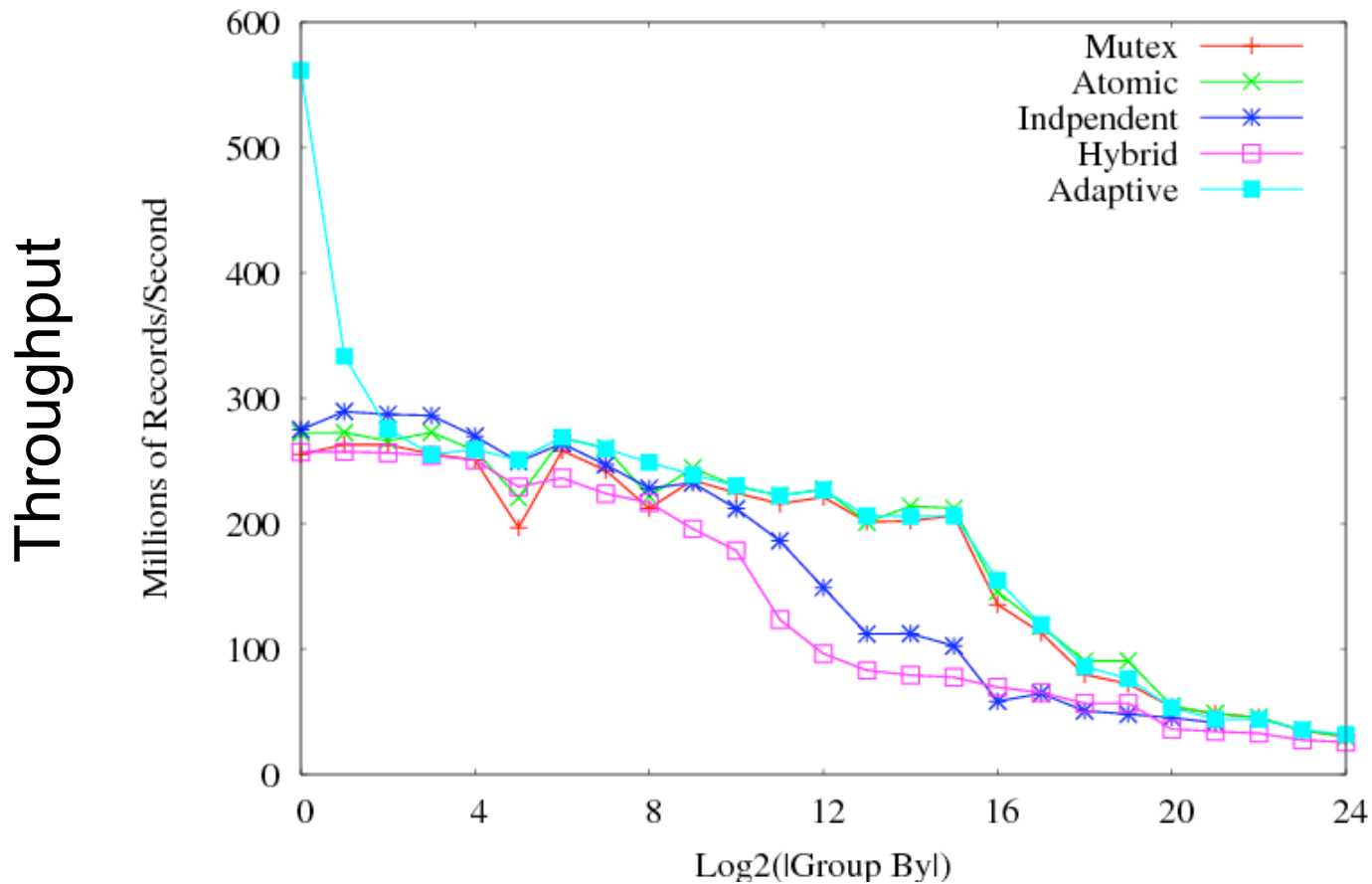
# Q1 (sum & count): Uniform Input

Adaptive Aggregation on Chip Multiprocessors

# Time Breakdown (Uniform Input)

Adaptive Aggregation on Chip Multiprocessors

# Q2 (min & max): Uniform Input

Adaptive Aggregation on Chip
Multiprocessors

# Q3 (DISTINCT): Uniform Input



Adaptive Aggregation on Chip Multiprocessors

# Q1 (sum & count): Self-similar



Throughput

Millions of Records/Second

Log2(|Group By|)

Legend:
- Mutex
- Atomic
- Independent
- Hybrid
- Adaptive

No Contention

Contention

Adaptive Aggregation on Chip
Multiprocessors

# See the paper for…

- The full contention model

- Results with other input distributions

- The impact of resampling the input stream
  - Able to adapt to changes in the distribution

- Scales with the number of computed aggregates
  - Global table with mutex / lock eventually out performs atomic instructions

# Conclusion

- Investigated aggregation performance on a real chip multiprocessor

- Identified locality and contention as key performance issues

- Introduced an adaptive aggregation operator that uses lightweight sampling to choose the best aggregation strategy

# Das Ende



Adaptive Aggregation on Chip Multiprocessors