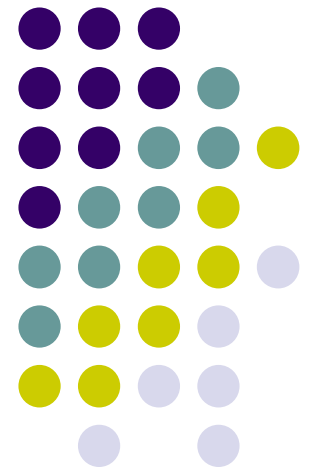




“Generic Database Cost Models for Hierarchical Memory Systems”, Manegold, Boncz, Kersten, VLDB’02

Cost models for Main-Memory database systems

Peter Boncz
Stefan Manegold
CWI (Amsterdam)





“Generic Database Cost Models for Hierarchical Memory Systems”, Manegold, Boncz, Kersten, VLDB’02



Contents

- Cost Components
- Disk-based vs. Main-Memory: Physical Costs
- Memory Access Pattern
- Conclusions





“Generic Database Cost Models for Hierarchical Memory Systems”, Manegold, Boncz, Kersten, VLDB’02



Introduction: Purposes of Cost Models

estimate / predict performance in order to

- understand and analyze algorithms
- tune algorithms
- optimize QEPs





Introduction: Cost Components

- **Logical Costs / Data Volume:**
 - estimate (intermediate) result sizes
 - “purely algebraic”, independent of algorithms
 - require statistics about the stored data
 - synopses, histograms, wavelets, sampling, etc.
- **Algorithmic Costs:**
 - complexity of algorithms
 - O -classes
- **Physical Costs:**
 - execution time, memory consumption, etc.
 - require knowledge about hardware





Physical Costs: Disk-based DBMS

- I/O costs (disk access) dominate execution
- CPU costs are negligible

$$\implies T = \#_{I/O} * t_{I/O}$$

- estimation of $\#_{I/O}$ is easy due to
- full control over buffer management:
 - known replacement strategy, pinning, exclusive access
 - full associativity \Rightarrow no conflicts





Physical Costs: Main-Memory DBMS

- no I/O
- only CPU costs and memory access costs
 - $\implies T = T_{CPU} + \#_{mem} * t_{mem}$
- estimation of $\#_{mem}$ is difficult due to
- no control over cache management:
 - unknown replacement strategy, no pinning, shared caches
 - limited associativity \Rightarrow conflicts
- analytical modeling of CPU costs is virtually impossible:
coding style, compiler optimization, runtime optimization (CPU), etc.





Memory Access Costs

- From Radix-work:
 - $T = T_{CPU} + M_{L1} * l_{L1} + M_{L2} * l_{L2} + M_{TLB} * l_{TLB}$
 - T_{CPU} is calibrated per operator
 - cache/memory characteristics are measured per system (The Calibrator)
- Problem:
 - hand-crafting the cost functions that estimate the cache miss rates is too expensive and not intuitive
 - we need a more generic solution





Memory Access Pattern: The Idea

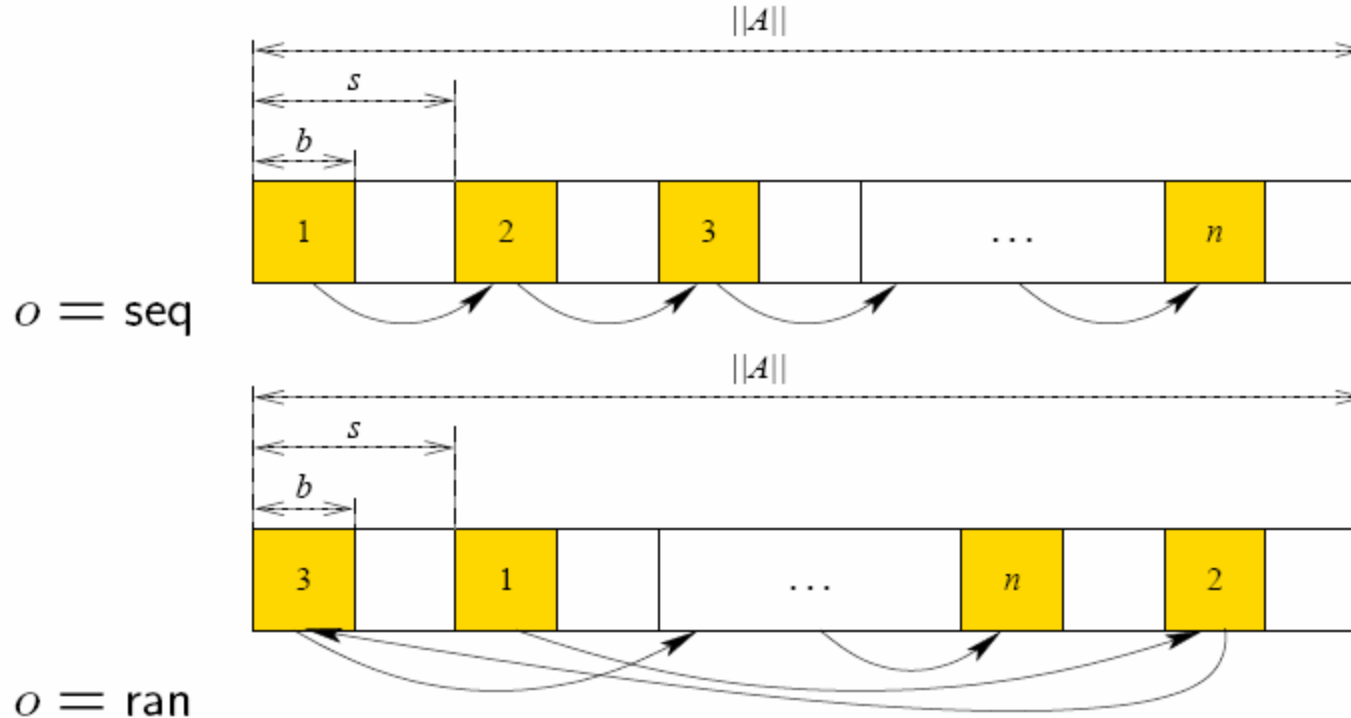
- specify basic memory access pattern and cost functions that calculate their cache miss rates
- use basic memory access pattern as building blocks to describe the memory access behavior of DB operators
⇒ compound memory access pattern
- provide rules to calculate the cache miss rates of compound memory access pattern using the cost functions of their building blocks





Simple Memory Access Pattern: $sMAP = (A, o, n, s, b,$

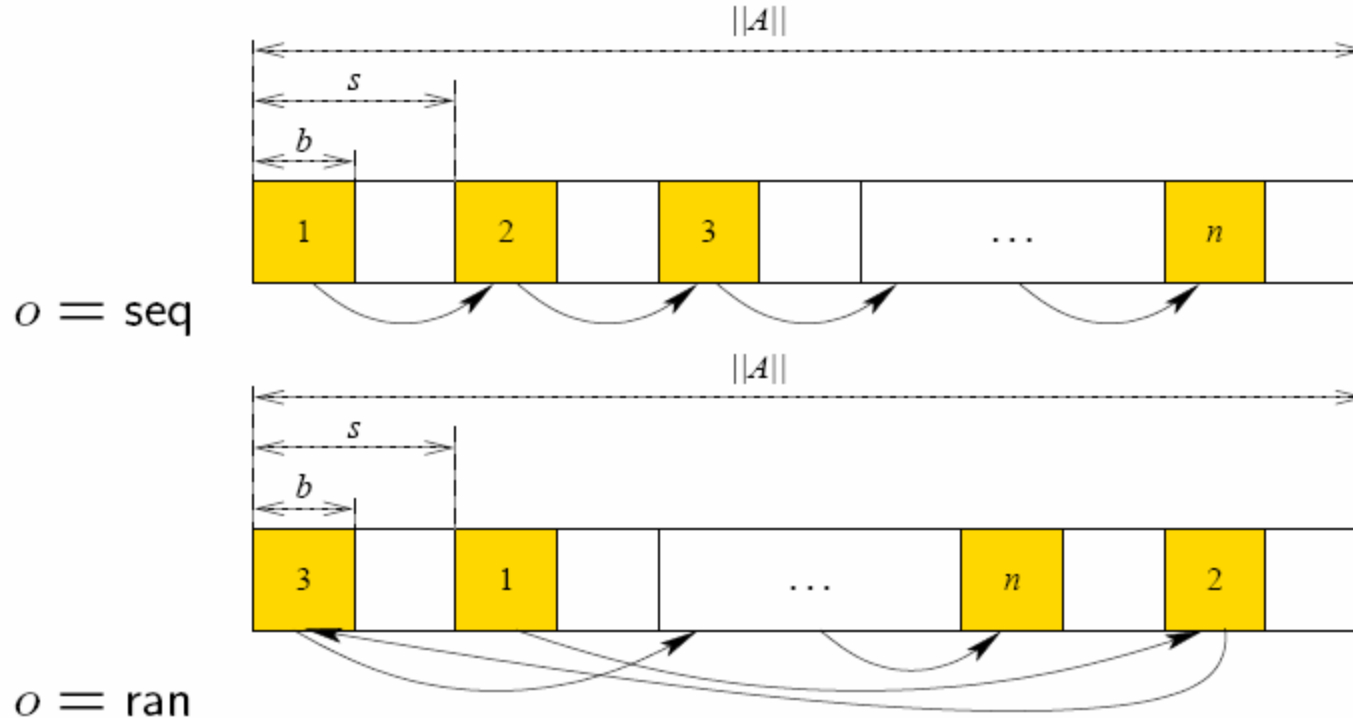
- traverse a memory area A once; $\|A\| = n \cdot s \wedge b \leq s$





Simple Memory Access Pattern: $sMAP = (A, o, n, s, b,$

- traverse a memory area A once $\|A\| = n \cdot s \wedge b \leq s$
 Region size





Simple Memory Access Pattern: Cache Misses

$$M_{L_i}((A, o, n, s, b)) =$$

$$\begin{cases} n \cdot \left(\left\lceil \frac{b}{LS_{L_i}} \right\rceil + \frac{(b-1) \bmod LS_{L_i}}{LS_{L_i}} \right), & \text{if } s - b \geq LS_{L_i} \\ |A|_{L_i}, & \text{if } s - b < LS_{L_i} \\ |A|_{L_i} + (n - \min\{|L_i|_{L_i}, |L_i|_s\}) \cdot \left(1 - \min\left\{1, \frac{\|L_i\|}{\|A\|}\right\} \right), & \text{if } s - b < LS_{L_i} \\ & \wedge o = \text{seq} \\ & \wedge o = \text{ran} \end{cases}$$

Skippy Access





Simple Memory Access Pattern: Cache Misses

$$M_{L_i}((A, o, n, s, b)) =$$

$$\left\{ \begin{array}{ll} n \cdot \left(\left\lceil \frac{b}{LS_{L_i}} \right\rceil + \frac{(b-1) \bmod LS_{L_i}}{LS_{L_i}} \right), & \text{if } s - b \geq LS_{L_i} \\ |A|_{L_i}, & \text{if } s - b < LS_{L_i} \\ & \wedge o = \text{seq} \\ |A|_{L_i} \\ + (n - \min\{|L_i|_{L_i}, |L_i|_s\}) \cdot \left(1 - \min\left\{1, \frac{\|L_i\|}{\|A\|}\right\} \right), & \text{if } s - b < LS_{L_i} \\ & \wedge o = \text{ran} \end{array} \right.$$

Full sequential access





Simple Memory Access Pattern: Cache Misses

$$M_{L_i}((A, o, n, s, b)) =$$

$$\left\{ \begin{array}{ll} n \cdot \left(\left\lceil \frac{b}{LS_{L_i}} \right\rceil + \frac{(b-1) \bmod LS_{L_i}}{LS_{L_i}} \right), & \text{if } s - b \geq LS_{L_i} \\ |A|_{L_i}, & \text{if } s - b < LS_{L_i} \\ & \wedge o = \text{seq} \\ |A|_{L_i} \quad \text{Random access: compulsory misses + data misses} \\ + (n - \min\{|L_i|_{L_i}, |L_i|_s\}) \cdot \left(1 - \min\left\{1, \frac{\|L_i\|}{\|A\|}\right\} \right), & \text{if } s - b < LS_{L_i} \\ & \wedge o = \text{ran} \end{array} \right.$$





Simple Memory Access Pattern: Cache Misses

Beware automatic prefetching!

Sequential access is much cheaper
-5GB/s vs 500MB/s
Workaround: use “amortized latency” for sequential accesses

$$M_{L_i}((A, o, n, s, b)) =$$

$$\left\{ \begin{array}{l} n \cdot \left(\left\lceil \frac{b}{LS_{L_i}} \right\rceil + \frac{(b-1) \bmod LS_{L_i}}{LS_{L_i}} \right), \\ |A|_{L_i}, \quad \text{Sequential memory access} \\ |A|_{L_i} \quad \text{Random memory access} \\ + (n - \min\{|L_i|_{L_i}, |L_i|_s\}) \cdot \left(1 - \min\left\{1, \frac{\|L_i\|}{\|A\|}\right\} \right), \end{array} \right. \begin{array}{l} \text{if } s - b < LS_{L_i} \\ \wedge o = \text{seq} \\ \text{if } s - b < LS_{L_i} \\ \wedge o = \text{ran} \end{array}$$





Repetitive Memory Access Pattern: $rMAP = (S, r, d)$

- repeat $S = (A, o, n, s, b)$ r -times

- for $o = \text{seq}$, we distinguish

$d = \text{uni}(-\text{directional})$ (“round-robin”)

$d = \text{bi}(-\text{directional})$ (“back-and-forth”)





Repetitive Memory Access Pattern: Cache Misses

$$M_{L_i}((S, r, d)) =$$

$$\begin{cases} M_{L_i}(S), & \text{if } M_{L_i}(S) \leq |L_i|_{L_i} \\ r \cdot M_{L_i}(S), & \text{if } M_{L_i}(S) > |L_i|_{L_i} \\ & \wedge d = \text{uni} \\ M_{L_i}(S) + (r - 1) \cdot (M_{L_i}(S) - K), & \text{else} \end{cases}$$

$$K = \begin{cases} |L_i|_{L_i}, & \text{if } o = \text{seq} \\ \frac{|L_i|_{L_i}}{M_{L_i}(S)} \cdot |L_i|_{L_i}, & \text{if } o = \text{ran} \end{cases}$$





Repetitive Memory Access Pattern: Cache Misses

$$M_{L_i}((S, r, d)) =$$

$$\left\{ \begin{array}{ll} M_{L_i}(S), & \text{Perfect reuse} \\ r \cdot M_{L_i}(S), & \text{if } M_{L_i}(S) > |L_i|_{L_i} \\ & \wedge d = \text{uni} \\ M_{L_i}(S) + (r - 1) \cdot (M_{L_i}(S) - K), & \text{else} \end{array} \right. \quad \text{if } M_{L_i}(S) \leq |L_i|_{L_i}$$

$$K = \begin{cases} |L_i|_{L_i}, & \text{if } o = \text{seq} \\ \frac{|L_i|_{L_i}}{M_{L_i}(S)} \cdot |L_i|_{L_i}, & \text{if } o = \text{ran} \end{cases}$$





Repetitive Memory Access Pattern: Cache Misses

$$M_{L^i}((S, r, d)) =$$

$$\begin{cases} M_{L^i}(S), & \text{if } M_{L^i}(S) \leq |L^i|_{L^i} \\ r \cdot M_{L^i}(S), & \text{No reuse} \quad \text{if } M_{L^i}(S) > |L^i|_{L^i} \\ & \quad \quad \quad \wedge d = \text{uni} \\ M_{L^i}(S) + (r - 1) \cdot (M_{L^i}(S) - K), & \text{else} \end{cases}$$

$$K = \begin{cases} |L^i|_{L^i}, & \text{if } o = \text{seq} \\ \frac{|L^i|_{L^i}}{M_{L^i}(S)} \cdot |L^i|_{L^i}, & \text{if } o = \text{ran} \end{cases}$$





Repetitive Memory Access Pattern: Cache Misses

$$M_{L_i}((S, r, d)) =$$

$$\begin{cases} M_{L_i}(S), & \text{if } M_{L_i}(S) \leq |L_i|_{L_i} \\ r \cdot M_{L_i}(S), & \text{if } M_{L_i}(S) > |L_i|_{L_i} \\ & \wedge d = \text{uni} \\ M_{L_i}(S) + (r - 1) \cdot (M_{L_i}(S) - K), & \text{else} \quad \text{“LRU” reuse} \end{cases}$$

$$K = \begin{cases} |L_i|_{L_i}, & \text{if } o = \text{seq} \\ \frac{|L_i|_{L_i}}{M_{L_i}(S)} \cdot |L_i|_{L_i}, & \text{if } o = \text{ran} \end{cases}$$





Compound Memory Access Pattern: Sequential Execution

- $\mathcal{P}_1, \mathcal{P}_2$ memory access pattern

$$M_{Li}(\mathcal{P}_1; \mathcal{P}_2) = M_{Li}(\mathcal{P}_1) + M_{Li}(\mathcal{P}_2)$$





Compound Memory Access Pattern: Concurrent Execution

- Q_1, Q_2 memory access pattern that do not revisit/re-use cache lines
- R_1, R_2 memory access pattern that do revisit/re-use cache lines

$$M_{Li}(Q_1||Q_2) = M_{Li}(Q_1) + M_{Li}(Q_2)$$

$$M_{Li}(Q_1||R_1) = M_{Li}(Q_1) + M_{Li}(R_1) = M_{Li}(R_1||Q_1)$$

$$M_{Li}(R_1||R_2) = M_{Li/2}(R_1) + M_{Li/2}(R_2)$$





Compound Memory Access Pattern: Concurrent Execution

- Q_1, Q_2 memory access pattern that do not revisit/re-use cache lines
- R_1, R_2 memory access pattern that do revisit/re-use cache lines

$$M_{Li}(Q_1||Q_2) = M_{Li}(Q_1) + M_{Li}(Q_2)$$

$$M_{Li}(Q_1||R_1) = M_{Li}(Q_1) + M_{Li}(R_1) = M_{Li}(R_1||Q_1)$$

$$M_{Li}(R_1||R_2) = M_{Li/2}(R_1) + M_{Li/2}(R_2)$$

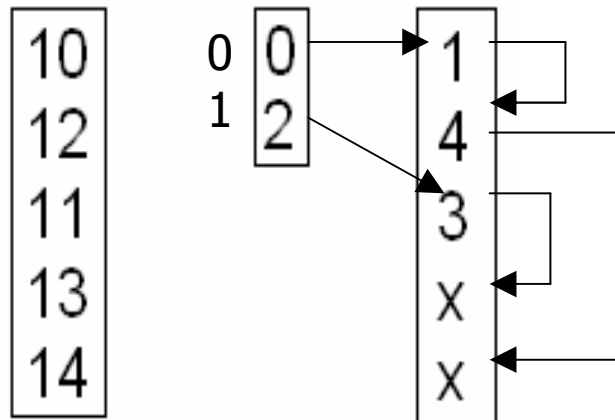




Compound Memory Access Pattern: Hash Join $B := L \bowtie R$

```
Build = ( (R, seq, R.BATcount, R.BUNsize, R.BUNsize)
          ||(R → link, seq, R.BATcount, 4, 4)
          ||(R → mask, ran, R.BATcount, 4, 4) )
```

keys buckets collision-list





Compound Memory Access Pattern: Hash Join $B := L \bowtie R$

$$\text{Build} = ((R, \text{seq}, R.\text{BATcount}, R.\text{BUNsize}, R.\text{BUNsize}) \\ || (R \rightarrow \text{link}, \text{seq}, R.\text{BATcount}, 4, 4) \\ || (R \rightarrow \text{mask}, \text{ran}, R.\text{BATcount}, 4, 4))$$

$$\text{Lookup} = ((R \rightarrow \text{mask}, \text{ran}, L.\text{BATcount}, 4, 4) \\ || (R \rightarrow \text{link}, \text{ran}, B.\text{BATcount}, 4, 4) \\ || (R, \text{ran}, B.\text{BATcount}, R.\text{BUNsize}, R.\text{BUNsize}))$$

$$\text{HashJoin} = (\text{Build}; ((L, \text{seq}, L.\text{BATcount}, L.\text{BUNsize}, L.\text{BUNsize}) \\ || \text{Lookup} \\ || (B, \text{seq}, B.\text{BATcount}, B.\text{BUNsize}, B.\text{BUNsize})))$$




Compound Memory Access Pattern: Hash Join $B := L \bowtie R$

Writes are non-blocking \perp Latency = 0

$$\text{Build} = ((R, \text{seq}, R.\text{BATcount}, R.\text{BUNsize}, R.\text{BUNsize}) \\ || (R \rightarrow \text{link}, \text{seq}, R.\text{BATcount}, 4, 4) \\ || (R \rightarrow \text{mask}, \text{ran}, R.\text{BATcount}, 4, 4))$$

$$\text{Lookup} = ((R \rightarrow \text{mask}, \text{ran}, L.\text{BATcount}, 4, 4) \\ || (R \rightarrow \text{link}, \text{ran}, B.\text{BATcount}, 4, 4) \\ || (R, \text{ran}, B.\text{BATcount}, R.\text{BUNsize}, R.\text{BUNsize}))$$

$$\text{HashJoin} = (\text{Build}; ((L, \text{seq}, L.\text{BATcount}, L.\text{BUNsize}, L.\text{BUNsize}) \\ || \text{Lookup} \\ || (B, \text{seq}, B.\text{BATcount}, B.\text{BUNsize}, B.\text{BUNsize})))$$

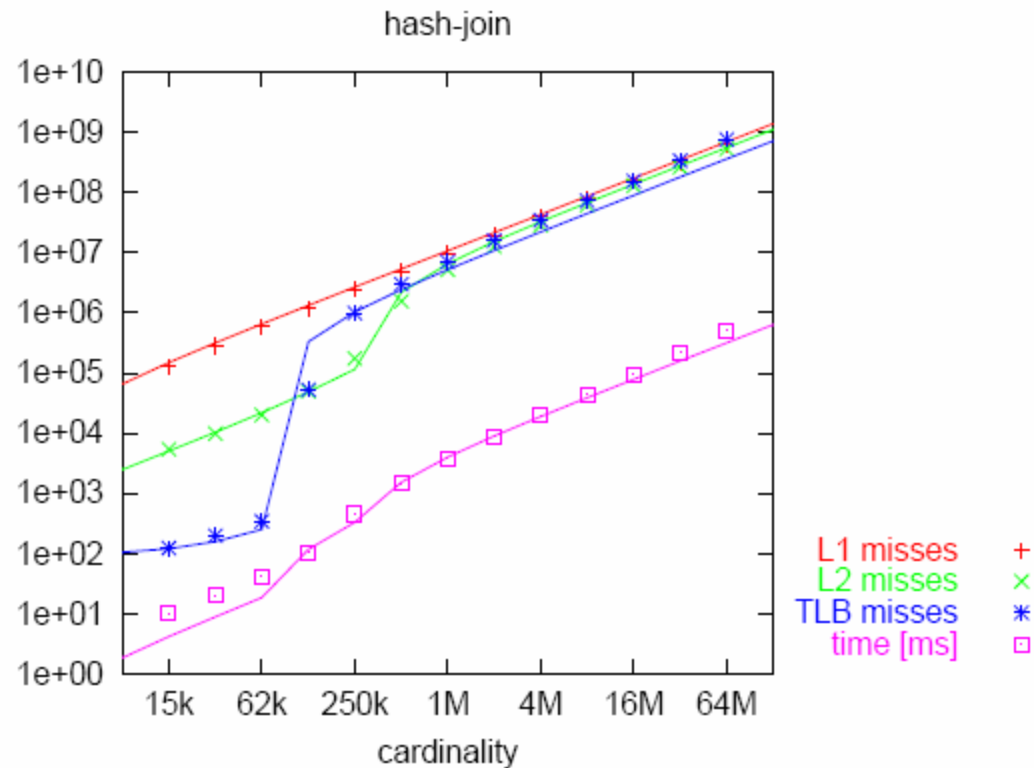



algorithm	pattern description	name
$W \leftarrow select(U)$	$s_tra(U) \odot s_tra(W)$	
$W \leftarrow nested_loop_join(U, V)$	$s_tra(U) \odot rs_tra(U.n, uni, V) \odot s_tra(W)$	$=:$ $nl_join(U, V, W)$
$H \leftarrow hash_build(V)$	$s_tra(V) \odot r_tra(H)$	$=:$ $build_hash(V, H)$
$W \leftarrow hash_probe(U, H)$	$s_tra(U) \odot r_acc(U.n, H) \odot s_tra(W)$	$=:$ $probe_hash(U, H, W)$
$W \leftarrow hash_join(U, V)$	$build_hash(V, H) \oplus probe_hash(U, H, W)$	$=:$ $h_join(U, V, W)$
$\{U_j\}_{j=1}^m \leftarrow cluster(U, m)$	$s_tra(U) \odot nest(\{U_j\}_{j=1}^m, m, s_tra(U_j), ran)$	$=:$ $part(U, m, \{U_j\}_{j=1}^m)$
$W \leftarrow part_nl_join(U, V, m)$	$part(U, m, \{U_j\}_{j=1}^m) \oplus part(V, m, \{V_j\}_{j=1}^m)$ $\oplus nl_join(U_1, V_1, W_1) \oplus \dots \oplus nl_join(U_m, V_m, W_m)$	
$W \leftarrow part_h_join(U, V, m)$	$part(U, m, \{U_j\}_{j=1}^m) \oplus part(V, m, \{V_j\}_{j=1}^m)$ $\oplus h_join(U_1, V_1, W_1) \oplus \dots \oplus h_join(U_m, V_m, W_m)$	



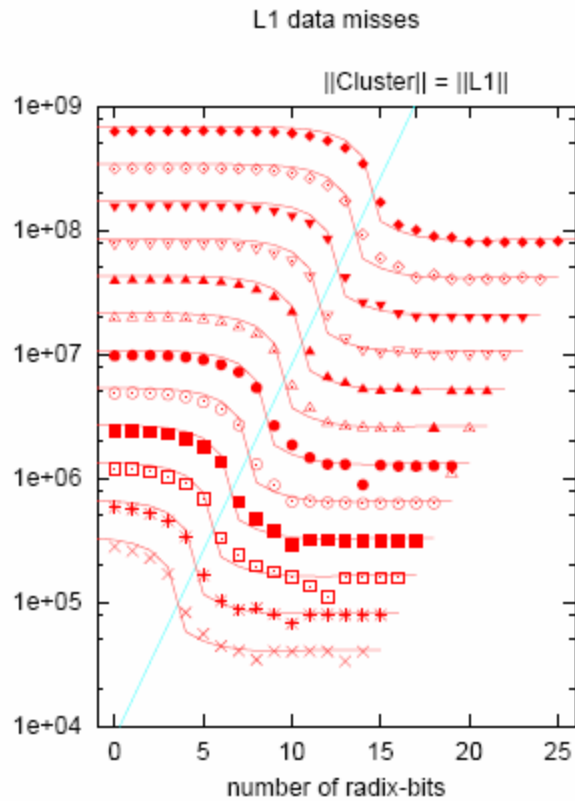


Compound Memory Access Pattern: Hash Join



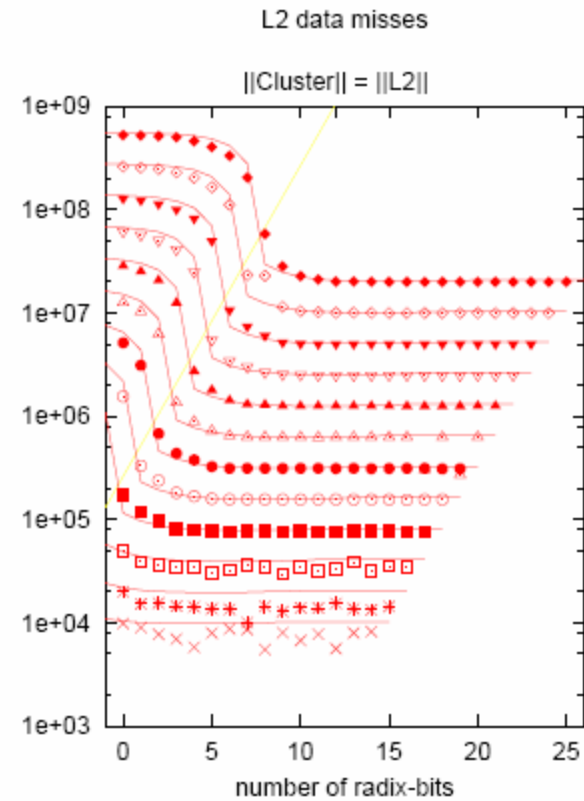


Compound Memory Access Pattern: Partitioned Hash Join



64000000 ◆
32000000 ◇

8000000 ▼
4000000 ▲



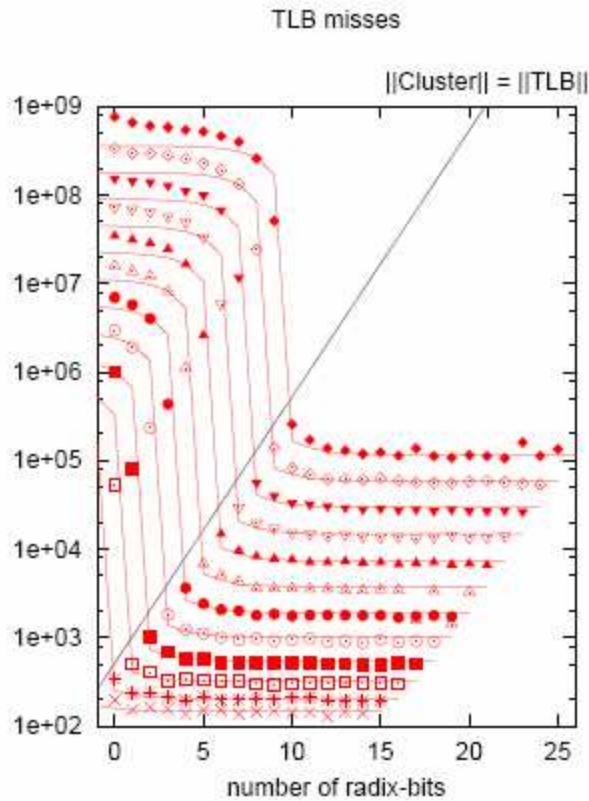
1000000 ●
500000 ○

125000 □
62500 *

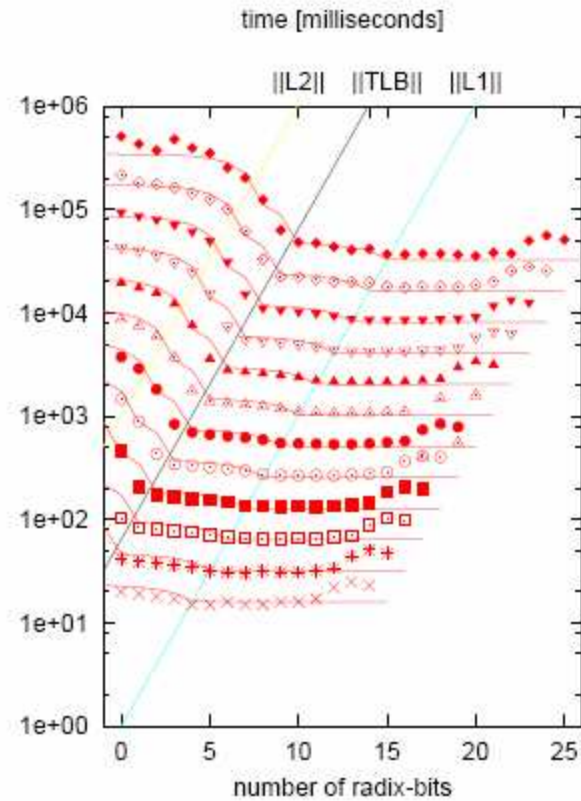




Compound Memory Access Pattern: Partitioned Hash Join



64000000	•	8000000	▽
32000000	◊	4000000	▲

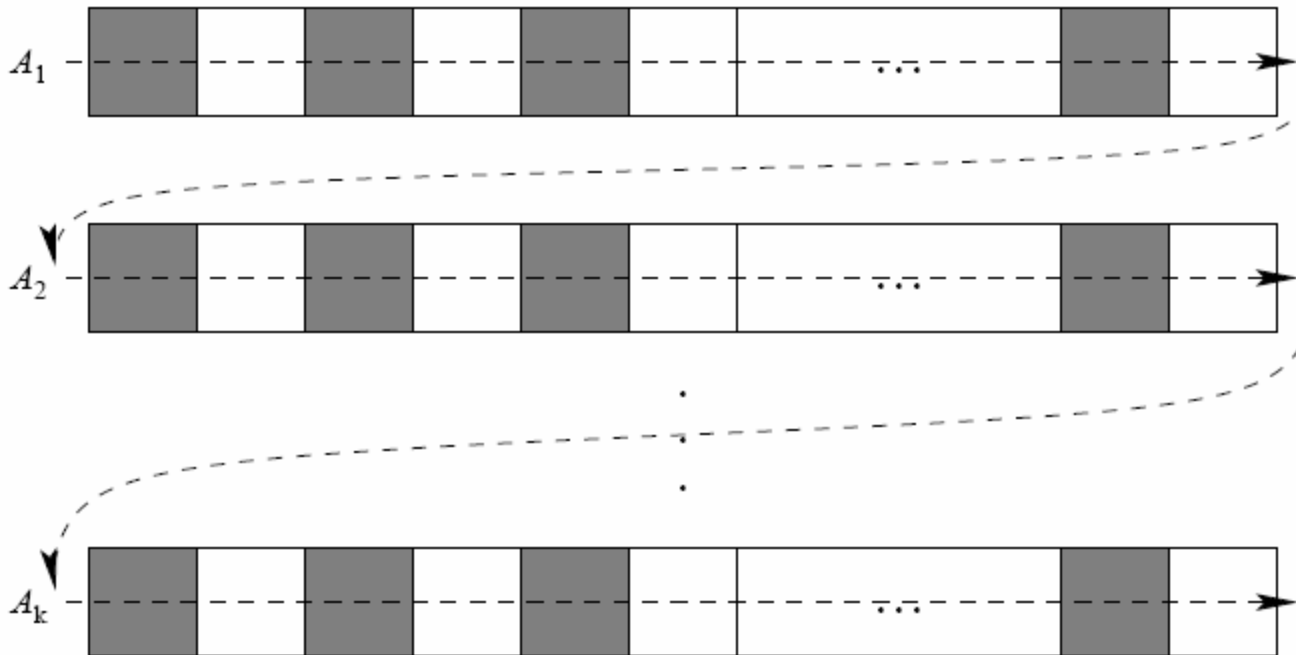


1000000	•	125000	◻
500000	◊	62500	*





Nested Memory Access Pattern: $nMAP = (\langle S_1, \dots, S_k \rangle, Y, O, D)$



$Y = \text{con}(\text{catenated}), \quad O = \text{seq}, \quad D = \text{uni}$





Nested Memory Access Pattern: Cache Misses

$$M_{Li}((\langle S_1, \dots, S_k \rangle, Y, O, D)) =$$

$$\begin{cases} k \cdot M_{Li}(S_j), & \text{if } Y = \text{con} \wedge s - b < LS_{Li} \\ M_{Li}(S(O)), & \text{if } Y = \text{int} \wedge o = \text{ran} \wedge n = 1 \\ M_{Li}(S(o)) + X & \text{if } Y = \text{int} \wedge o = \text{seq} \wedge \\ & s - b < LS_{Li} \wedge k \cdot \left\lceil \frac{b}{LS_{Li}} \right\rceil > |Li|_{Li} \\ M_{Li}(S(o)), & \text{else} \end{cases}$$

$$S(x) = \left(\bigcup_{j=1}^k A_j, x, k \cdot n, s, b \right)$$

$$X = (n - 1) \cdot (k - h), \quad h = \begin{cases} 0, & \text{if } O = \text{seq} \wedge D = \text{uni} \\ |TLB|, & \text{if } O = \text{seq} \wedge D = \text{bi} \\ \frac{|TLB|}{k} \cdot |TLB|, & \text{if } O = \text{ran} \end{cases}$$





Compound Memory Access Pattern: Radix Cluster

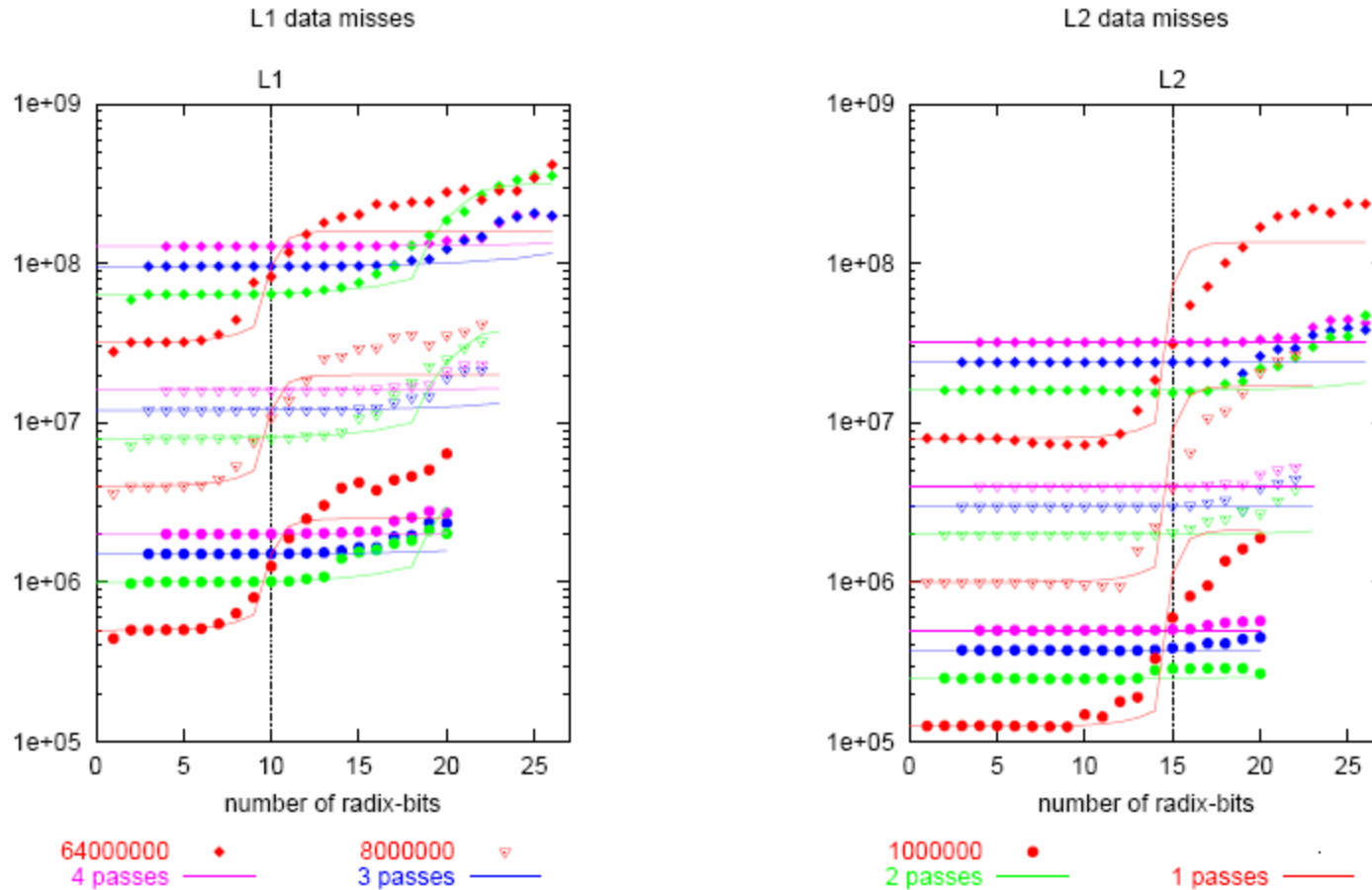
$$S_j = (S_j, \text{seq}, S_j.\text{BATcount}, S_j.\text{BUNsize}, S_j.\text{BUNsize})$$

$$\text{RadixCluster} = ((R, \text{seq}, R.\text{BATcount}, R.\text{BUNsize}, R.\text{BUNsize}) \\ || (< S_1, \dots, S_k >, \text{int}, \text{ran}))$$





Compound Memory Access Pattern: Radix Cluster

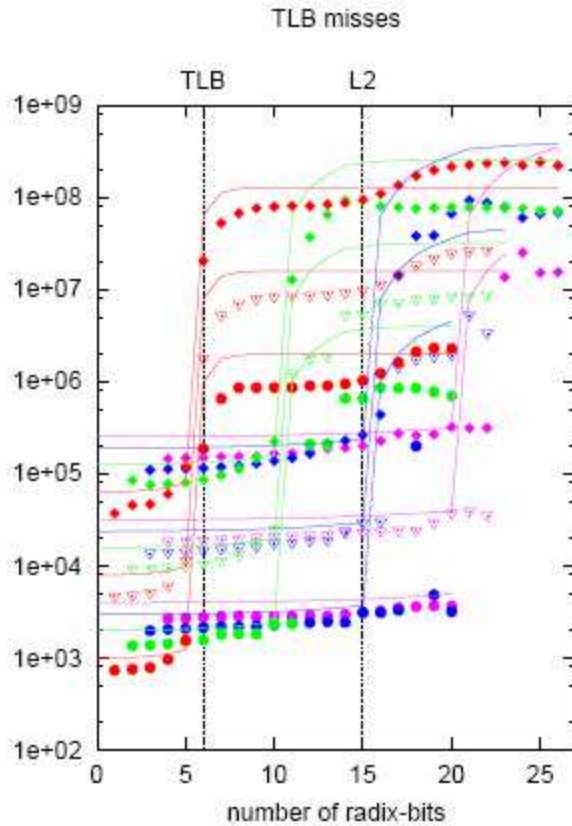




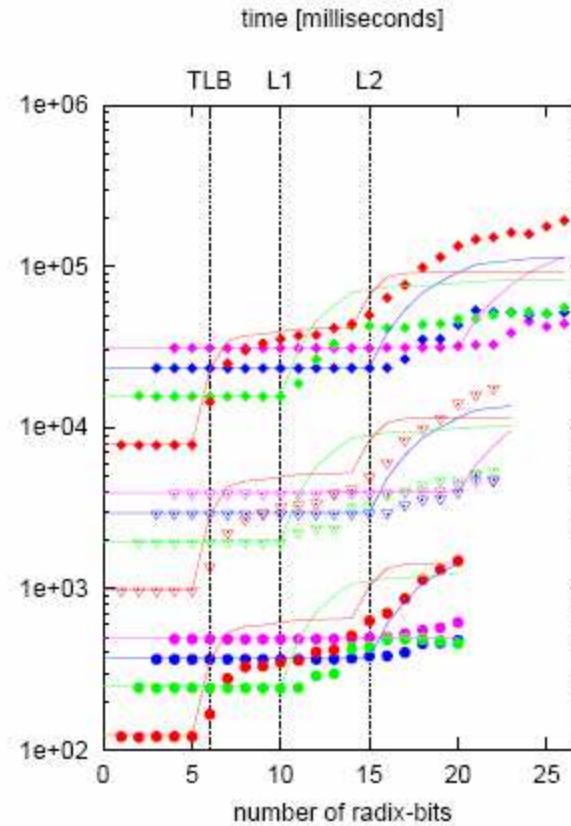
“Generic Database Cost Models for Hierarchical Memory Systems”, Manegold, Boncz, Kersten, VLDB’02



Compound Memory Access Pattern: Radix Cluster



64000000 4 passes (red diamond)
8000000 3 passes (blue inverted triangle)



1000000 2 passes (green circle)
1 passes (red circle)





Conclusions

- generic, architecture-independent main-memory cost-modeling is possible (at least in Monet)
- cache impact is crucial for intra-operator costs
- use operator cost functions to always pick *the best* algorithm
⇒ optimization on higher levels can use simpler cost models

