# CS5226 Lecture 4
# Memory Tuning

# Memory Management
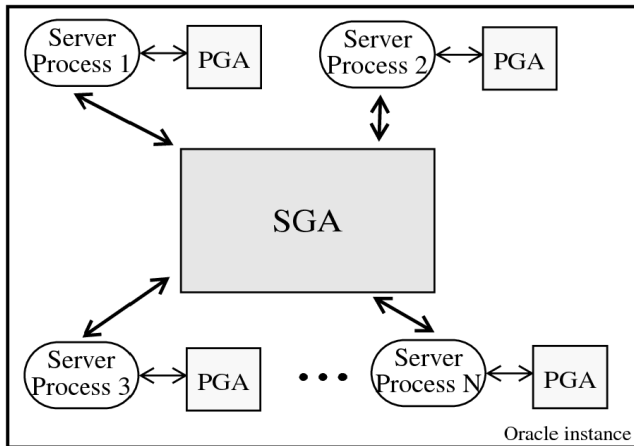
- **Shared data and control structures**:
    - buffer cache
    - catalog metadata
    - log buffer
    - lock structures
    - SQL execution plans
    - etc
- **Local data and control structures**:
    - operator work area = memory allocated for evaluating an SQL operator
    - Memory-intensive operators:
        - ★ sort
        - ★ hash join

# Oracle Memory Model



(Dageville & Zait, VLDB 2002)

System Global Area (SGA) = shared memory region
Process Global Area (PGA) = private memory region

# Oracle Memory Management

- **Pre 9i** - Manual memory mangement

  - **Parameters for tuning SGA**:
    - ★ DB_CACHE_SIZE, SHARED_POOL_SIZE, LARGE_POOL_SIZE, JAVA_POOL_SIZE, etc.

  - **Parameters for tuning PGA**:
    - ★ SORT_AREA_SIZE, HASH_AREA_SIZE, BITMAP_MERGE_AREA_SIZE, etc.

- **9i** - Automatic PGA memory management

  - PGA_AGGREGATE_TARGET parameter

- **10g** - Automatic SGA memory management

  - SGA_TARGET parameter

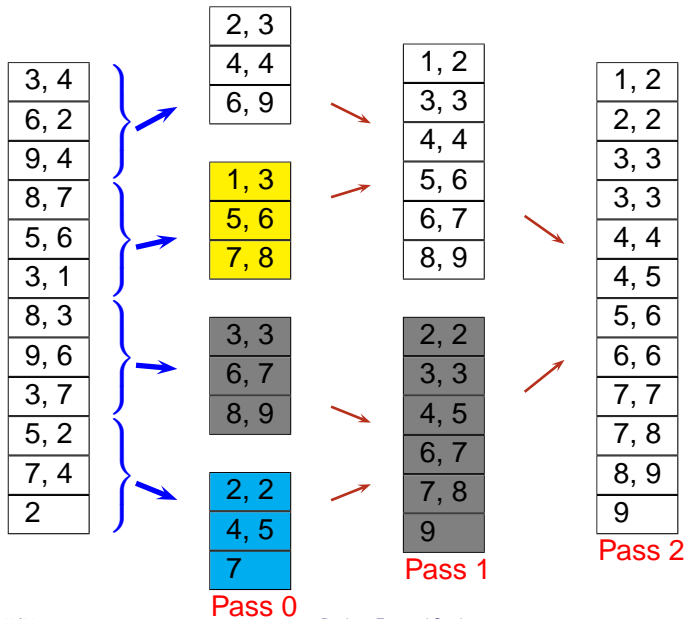- **11g** - Automatic memory management

  - MEMORY_TARGET parameter

# Memory-Intensive Operators

- Examples of memory-intensive operators:
  - sort
  - hash join
- How to manage memory for such operators?

# Review: External Merge Sort

- Sorting file of *N* pages with B buffer pages

- **Pass 0**: Creation of sorted runs

  - Read in and sort B pages at a time
  - Number of sorted runs created = $\lceil N/B \rceil$
  - Size of each sorted run = *B* pages (except possibly for last run)

- **Pass i,** $i \geq 1$: Merging of sorted runs

  - Use $B - 1$ buffer pages for input & one buffer page for output
  - Performs (B-1)-way merge

- Analysis:

  - $N_0$ = number of initial sorted runs = $\lceil N/B \rceil$
  - total number of passes = $\lceil \log_{B-1}(N_0) \rceil + 1$

# External Merge Sort with B=3: Example



| | | | | | | |
|---|---|---|---|---|---|---|
| 3, 4 | | 2, 3 | | 1, 2 | | 1, 2 |
| 6, 2 | | 4, 4 | | 3, 3 | | 2, 2 |
| 9, 4 | | 6, 9 | | 4, 4 | | 3, 3 |
| 8, 7 | | 1, 3 | | 5, 6 | | 3, 3 |
| 5, 6 | | 5, 6 | | 6, 7 | | 4, 4 |
| 3, 1 | | 7, 8 | | 8, 9 | | 4, 5 |
| 8, 3 | | 3, 3 | | 2, 2 | | 5, 6 |
| 9, 6 | | 6, 7 | | 3, 3 | | 6, 6 |
| 3, 7 | | 8, 9 | | 4, 5 | | 7, 7 |
| 5, 2 | | 2, 2 | | 6, 7 | | 7, 8 |
| 7, 4 | | 4, 5 | | 7, 8 | | 8, 9 |
| 2 | | 7 | | 9 | | 9 |

Pass 0   Pass 1   Pass 2

# Review: Hash Join, $R \bowtie S$

- ▶ Idea:
  - ▶ Partition $R$ and $S$ into $k$ partitions using some hash function $h$
    - ★ $R = R_1 \cup R_2 \cup \cdots \cup R_k$
    - ★ $S = S_1 \cup S_2 \cup \cdots \cup S_k$
  - ▶ Joins corresponding pair of partitions
    - ★ $R \bowtie S = (R_1 \bowtie S_1) \cup (R_2 \bowtie S_2) \cup \cdots \cup (R_k \bowtie S_k)$
- ▶ Algorithms:
  - ▶ Grace hash join
  - ▶ Hybrid hash join

# Review: Grace Hash Join $R \bowtie S$

► Consists of three phases:

1. Partition $R$ into $R_1, \cdots, R_k$
2. Partition $S$ into $S_1, \cdots, S_k$
3. Probing phase: probes each $R_i$ with $S_i$
   * Read $R_i$ to build a hash table
   * Read $S_i$ to probe hash table

► R is called the build relation & S is called the probe relation

Partitioning (building) phases
initialize a hash table $T$ with $k$ buckets
for each tuple $r \in R$ do
    insert $r$ into bucket $h(r)$ of $T$
write each bucket $R_i$ of $T$ to disk
initialize a hash table $T$ with $k$ buckets
for each tuple $s \in S$ do
    insert $s$ into bucket $h(s)$ of $T$
write each bucket $S_i$ of $T$ to disk

Probing (matching) phase

for i = 1 to k do
    initialize a hash table $T$
    for each tuple $r$ in partition $R_i$ do
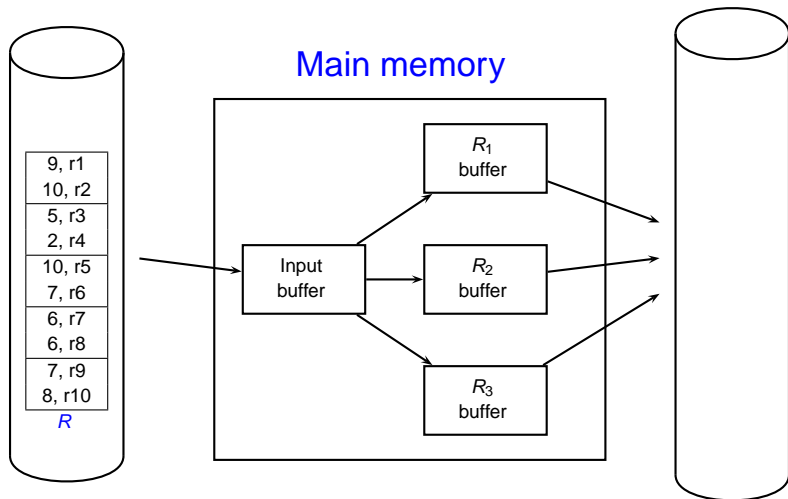        insert $r$ into bucket $h'(r)$ of $T$
    for each tuple $s$ in partition $S_i$ do
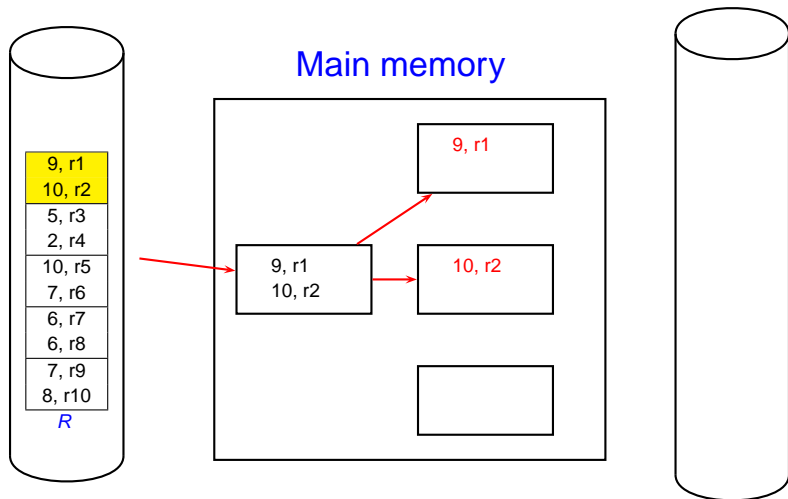        for each tuple $r$ in bucket $h'(s)$ of $T$ do
            if $r$ and $s$ matches then output $(r, s)$
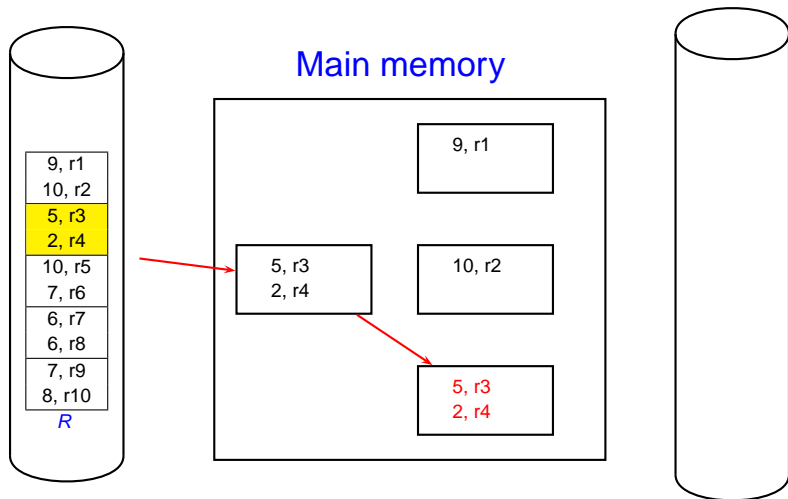
# Grace Hash Join: Partitioning Relation R



Main memory

| 9, r1 |
| 10, r2 |
| 5, r3 |
| 2, r4 |
| 10, r5 |
| 7, r6 |
| 6, r7 |
| 6, r8 |
| 7, r9 |
| 8, r10 |

$R$

Input buffer

$R_1$ buffer

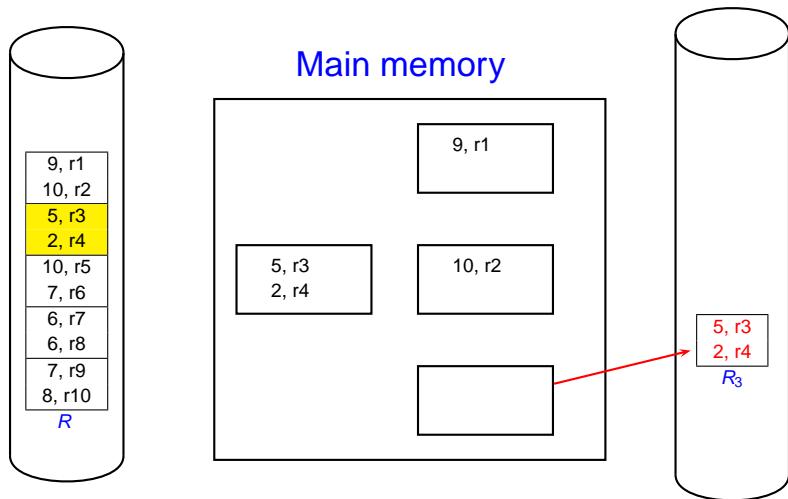$R_2$ buffer

$R_3$ buffer

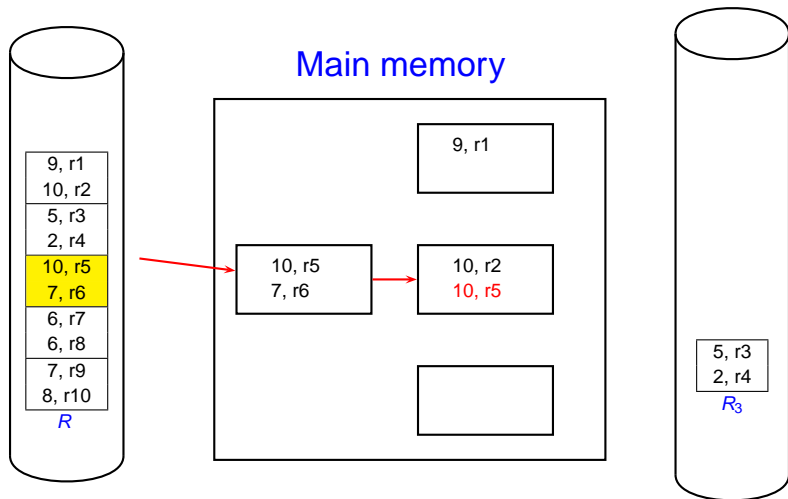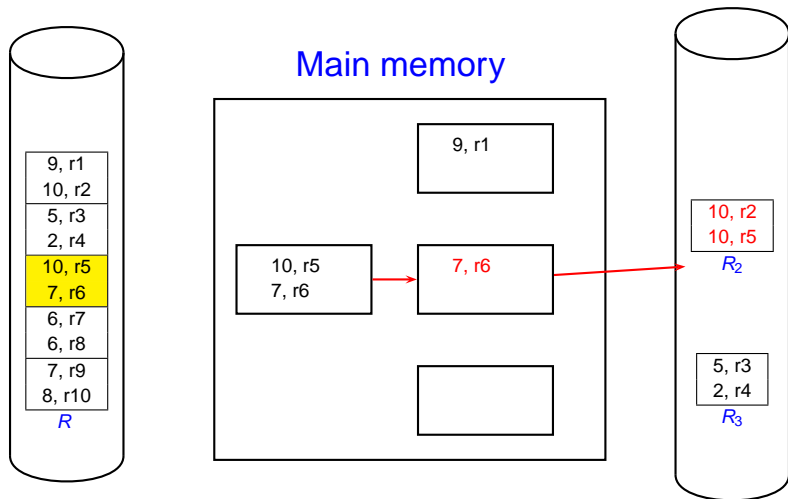Hash function: $h(v) = v \mod 3$

# Grace Hash Join: Partitioning Relation R

# Grace Hash Join: Partitioning Relation R

# Grace Hash Join: Partitioning Relation R

# Grace Hash Join: Partitioning Relation R



Main memory

| 9, r1 |
| 10, r2 |
| 5, r3 |
| 2, r4 |
| 10, r5 |
| 7, r6 |
| 6, r7 |
| 6, r8 |
| 7, r9 |
| 8, r10 |
*R*

9, r1

10, r5
7, r6

10, r2
10, r5

5, r3
2, r4
*R₃*
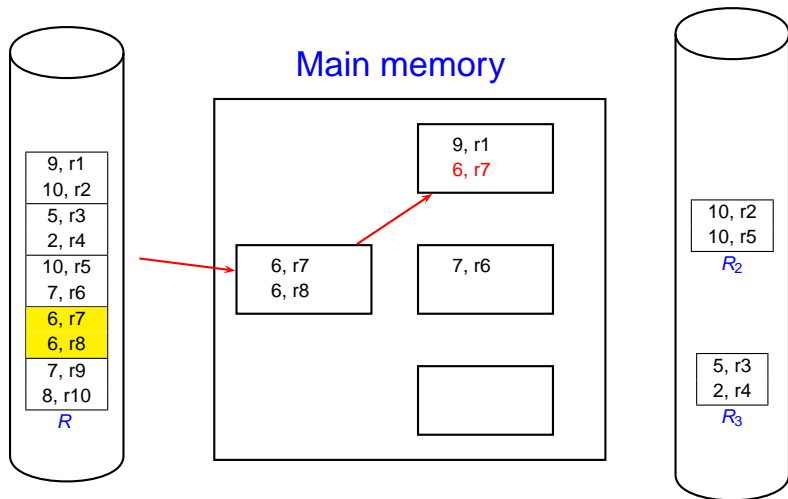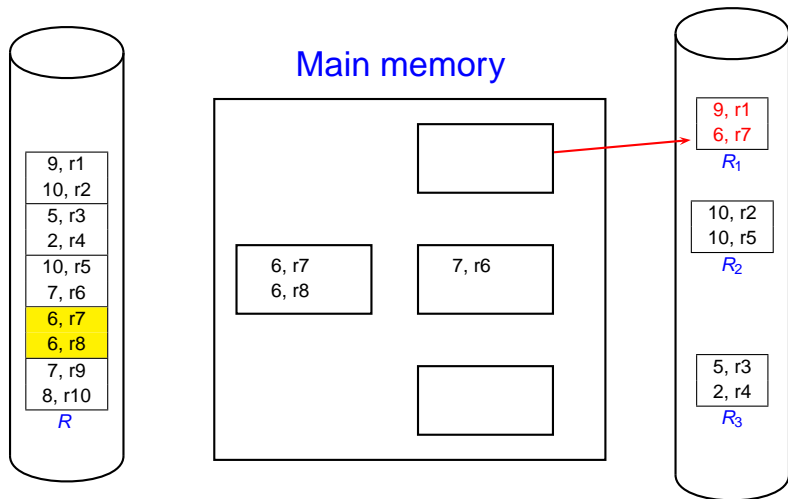
# Grace Hash Join: Partitioning Relation R

Review: Grace Hash Join

10

# Grace Hash Join: Partitioning Relation R



Main memory

# Grace Hash Join: Partitioning Relation R

Main memory

| 9, r1 |
| 10, r2 |
| 5, r3 |
| 2, r4 |
| 10, r5 |
| 7, r6 |
| 6, r7 |
| 6, r8 |
| 7, r9 |
| 8, r10 |
| *R* |

| 6, r7 |
| 6, r8 |

| 7, r6 |

| 9, r1 |
| 6, r7 |
| *R₁* |

| 10, r2 |
| 10, r5 |
| *R₂* |

| 5, r3 |
| 2, r4 |
| *R₃* |

# Grace Hash Join: Partitioning Relation R



Main memory

| | |
|---|---|
| 9, r1 | |
| 10, r2 | |
| 5, r3 | |
| 2, r4 | |
| 10, r5 | 7, r9 |
| 7, r6 | 8, r10 |
| 6, r7 | |
| 6, r8 | |
| 7, r9 | |
| 8, r10 | |
| R | |

7, r6
7, r9

8, r10

9, r1
6, r7
$R_1$

10, r2
10, r5
$R_2$

5, r3
2, r4
$R_3$

# Grace Hash Join: Partitioning Relation R



Main memory

| R |
|---|
| 9, r1 |
| 10, r2 |
| 5, r3 |
| 2, r4 |
| 10, r5 |
| 7, r6 |
| 6, r7 |
| 6, r8 |
| 7, r9 |
| 8, r10 |

7, r9
8, r10

| $R_1$ |
|---|
| 9, r1 |
| 6, r7 |

| $R_2$ |
|---|
| 10, r2 |
| 10, r5 |
| 7, r6 |
| 7, r9 |

| $R_3$ |
|---|
| 5, r3 |
| 2, r4 |
| 8, r10 |

*R* is partitioned into $R_1$, $R_2$, & $R_3$

# Grace Hash Join: Partitioning Relation S



Main memory

S₁ buffer

Input buffer

S₂ buffer

S₃ buffer

3, s4
18, s6
9, s10
$S_1$

10, s2
7, s3
10, s7
10, s8
4, s9
$S_2$

5, s1
2, s5
$S_3$

5, s1
10, s2
7, s3
3, s4
2, s5
18, s6
10, s7
10, s8
4, s9
9, s10
$S$

Similarly, $S$ is partitioned into $S_1$, $S_2$, & $S_3$

# Grace Hash Join: Probing Phase

Main memory

| 9, r1 |
| 6, r7 |
| $R_1$ |

| 3, s4 |
| 18, s6 |
| 9, s10 |
| $S_1$ |

| 10, r2 |
| 10, r5 |
| 7, r6 |
| 7, r9 |
| $R_2$ |

| 10, s2 |
| 7, s3 |
| 10, s7 |
| 10, s8 |
| 4, s9 |
| $S_2$ |

| 5, r3 |
| 2, r4 |
| 8, r10 |
| $R_3$ |

| 5, s1 |
| 2, s5 |
| $S_3$ |

Input buffer

Output buffer

Hash table

# Grace Hash Join: Probing Phase



Main memory

| 9, r1 |
| 6, r7 |

build

Hash table

| (9,r1) |
| (6,r7) |
|  |

$R_1$: 9,r1 / 6,r7

$S_1$: 3, s4 / 18, s6 / 9, s10

$R_2$: 10, r2 / 10, r5 / 7, r6 / 7, r9

$S_2$: 10, s2 / 7, s3 / 10, s7 / 10, s8 / 4, s9

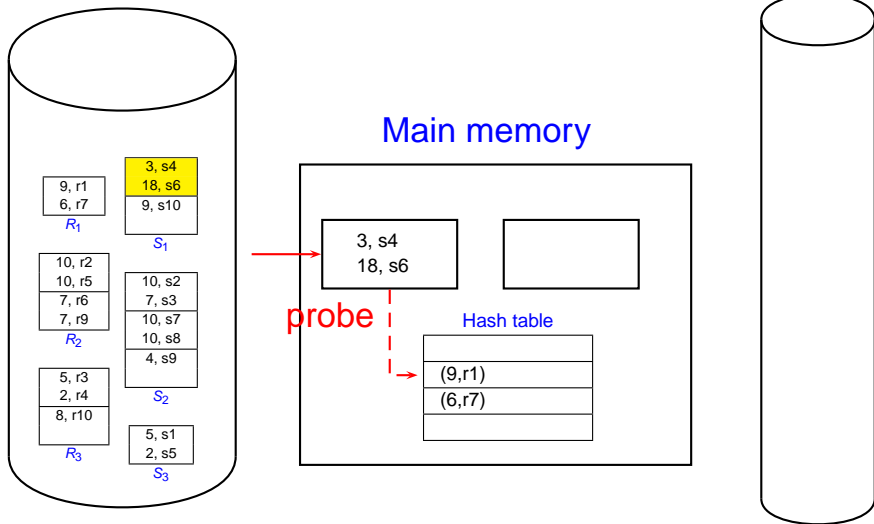$R_3$: 5, r3 / 2, r4 / 8, r10

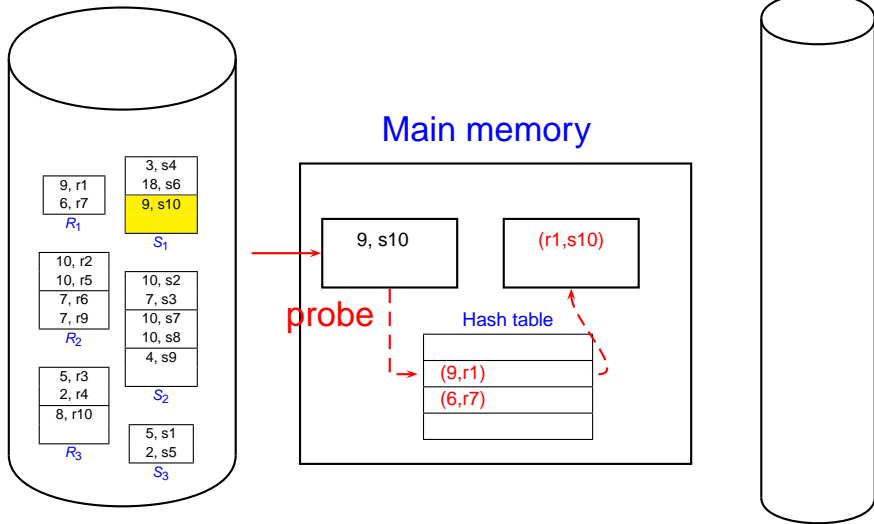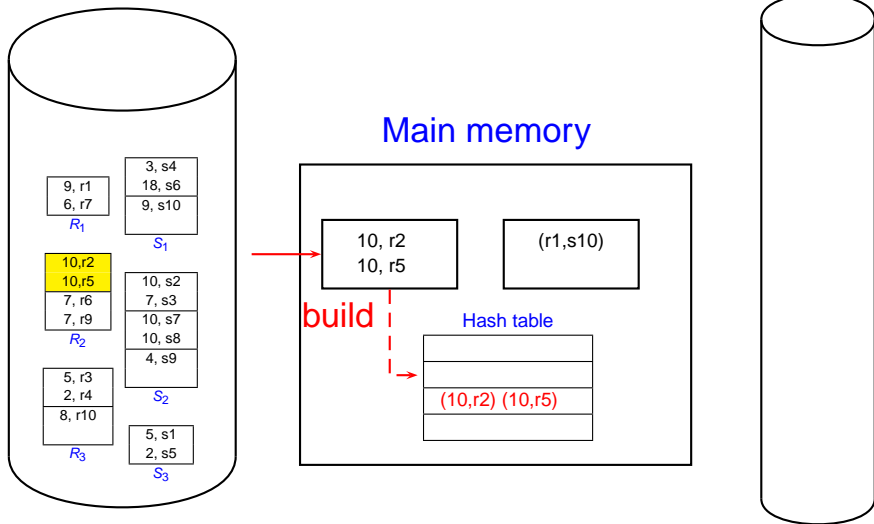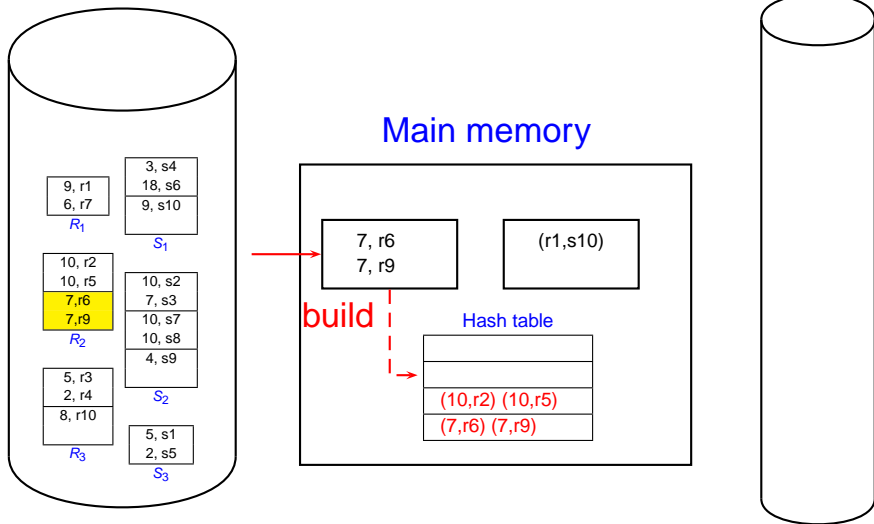$S_3$: 5, s1 / 2, s5

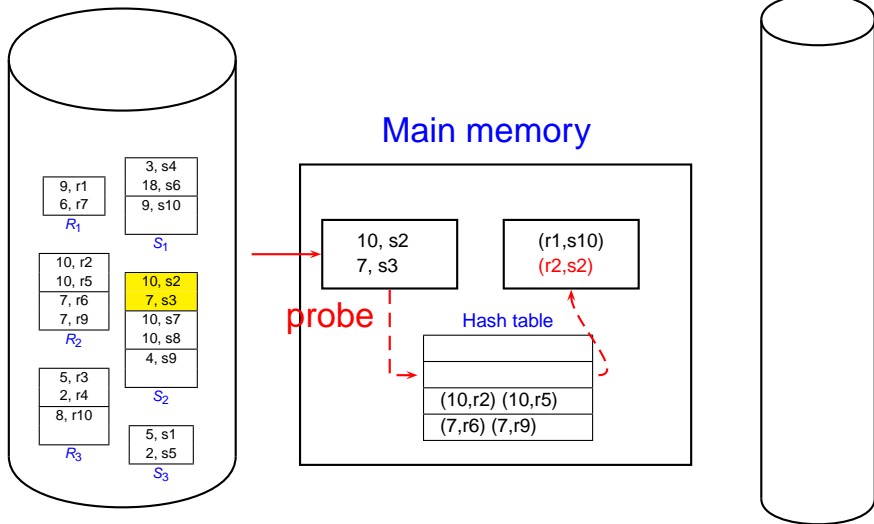Hash function: $h'(v) = v \bmod 4$

# Grace Hash Join: Probing Phase

# Grace Hash Join: Probing Phase

# Grace Hash Join: Probing Phase



Main memory

| 9, r1 | | 3, s4 |
| 6, r7 | | 18, s6 |
| $R_1$ | | 9, s10 |
| | | $S_1$ |

10, r2
10, r5

(r1,s10)

build

Hash table

(10,r2) (10,r5)

| 10,r2 | | 10, s2 |
| 10,r5 | | 7, s3 |
| 7, r6 | | 10, s7 |
| 7, r9 | | 10, s8 |
| $R_2$ | | 4, s9 |
| | | $S_2$ |

| 5, r3 | | |
| 2, r4 | | 5, s1 |
| 8, r10 | | 2, s5 |
| $R_3$ | | $S_3$ |

# Grace Hash Join: Probing Phase



Main memory

| 9, r1 |
| 6, r7 |
| $R_1$ |

| 3, s4 |
| 18, s6 |
| 9, s10 |
| $S_1$ |

| 10, r2 |
| 10, r5 |
| 7,r6 |
| 7,r9 |
| $R_2$ |

| 10, s2 |
| 7, s3 |
| 10, s7 |
| 10, s8 |
| 4, s9 |
| $S_2$ |

| 5, r3 |
| 2, r4 |
| 8, r10 |
| $R_3$ |

| 5, s1 |
| 2, s5 |
| $S_3$ |

| 7, r6 |
| 7, r9 |

(r1,s10)

build

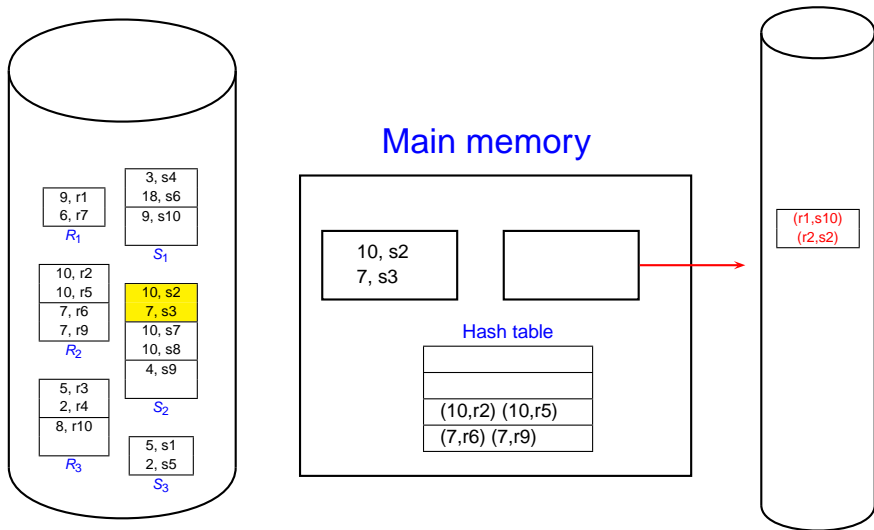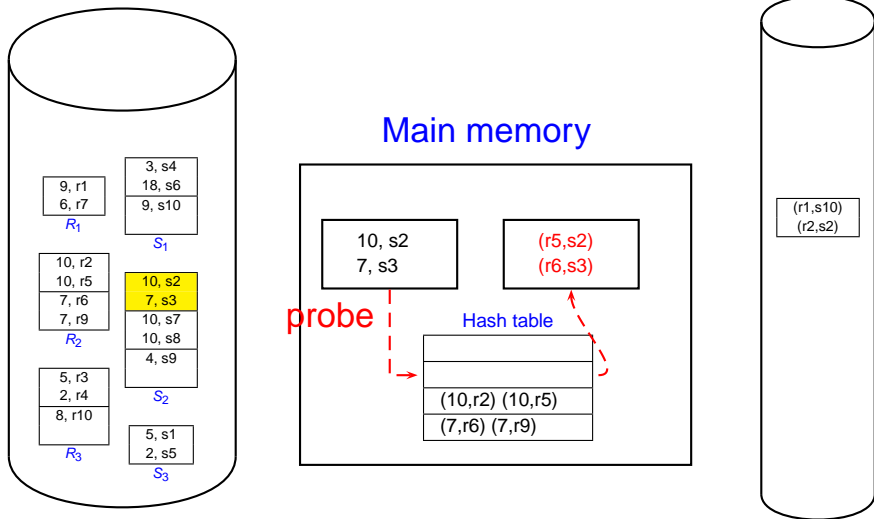Hash table

| (10,r2) (10,r5) |
| (7,r6) (7,r9) |

# Grace Hash Join: Probing Phase

# Grace Hash Join: Probing Phase

# Grace Hash Join: Probing Phase



Main memory

probe

Hash table

(r5,s2)
(r6,s3)

10, s2
7, s3

(10,r2) (10,r5)
(7,r6) (7,r9)

9, r1
6, r7
$R_1$

3, s4
18, s6
9, s10
$S_1$

10, r2
10, r5
7, r6
7, r9
$R_2$

10, s2
7, s3
10, s7
10, s8
4, s9
$S_2$

5, r3
2, r4
8, r10
$R_3$

5, s1
2, s5
$S_3$

(r1,s10)
(r2,s2)

# Grace Hash Join: Probing Phase



Main memory

Hash table

# Grace Hash Join: Probing Phase

# Grace Hash Join: Probing Phase



Main memory

| | |
|---|---|
| 9, r1 | 3, s4 |
| 6, r7 | 18, s6 |
| $R_1$ | 9, s10 |
| | $S_1$ |

| 10, r2 | 10, s2 |
| 10, r5 | 7, s3 |
| 7, r6 | 10, s7 |
| 7, r9 | 10, s8 |
| $R_2$ | 4, s9 |
| | $S_2$ |

| 5, r3 | 5, s1 |
| 2, r4 | 2, s5 |
| 8, r10 | $S_3$ |
| $R_3$ | |

10, s2
7, s3

(r9,s3)

probe

Hash table

| (10,r2) (10,r5) |
| (7,r6) (7,r9) |

| (r1,s10) |
| (r2,s2) |
| (r5,s2) |
| (r6,s3) |

Continue with probe $R_2$, build $R_3$, & probe $R_3$

# Review: Hybrid Hash Join (HHJ), $R \bowtie S$

- ▶ Improvement of Grace Hash Join

- ▶ Enables some build partitions $R_i$ to be resident in memory at the end of partitioning build relation $R$

  - ▶ Each $R_i$ is either a resident or spilled partition
  - ▶ Memory usage
    - ⋆ hash tables for resident partitions
    - ⋆ output buffers for spilled partitions
    - ⋆ input buffer
    - ⋆ free buffer - spool area for pages flushed to disk

- ▶ During partitioning of probe table $S$, $S_i$ can join with $R_i$ if $R_i$ is resident partition

- ▶ Number of partitions, k = $\sqrt{F|R|}$

  - ▶ $|R|$ = size of build relation R in pages
  - ▶ F = fudge factor

# HHJ: Partition Build Relation

01. let k = $\sqrt{F|R|}$
02. for i = 1 to *k* do
03.     mark $R_i$ as a resident partition & allocate one page to $R_i$
04. allocate input buffer & assign remaining pages to free buffer
05. for each tuple $t \in R$ do
06.     let *t* be hashed to $R_i$
07.     if ($R_i$ is full) and ($R_i$ is resident) then
08.         if (there's a free page) then
09.             add free page to $R_i$
10.         else
11.             mark $R_i$ as spilled
12.             flush $R_i$ to disk & free all pages (except one) of $R_i$
13.     else if ($R_i$ is full) and ($R_i$ is spilled) then
14.         flush $R_i$ to disk
15.     insert *t* into $R_i$
16. flush & free all pages in spilled partitions

# HHJ: Partition Probe Relation

01. for i = 1 to $k$ do
02.   if ($R_i$ is spilled) then
03.     allocate one page to $S_i$
04. for each tuple $t \in S$ do
05.   let $t$ be hashed to $S_i$
06.   if ($R_i$ is resident) then
07.     use $t$ to probe $R_i$ for matches
08.   else
09.     insert $t$ into $S_i$
10.     if ($S_i$ is full) then
11.       flush $S_i$ to disk
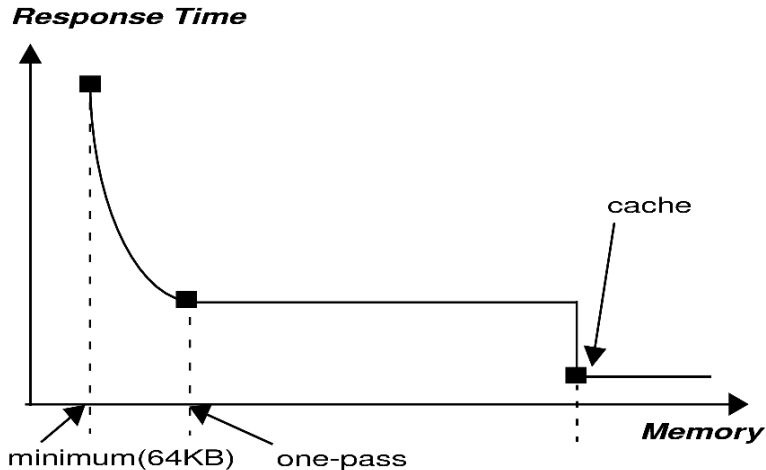12. flush all pages in $S$ partitions & free all pages

# HHJ: Probing Phase

01. for i = 1 to *k* do
02.    if (*R$_i$* is spilled) then
03.       read in *R$_i$* & build a hash table
04.       for each tuple *t* ∈ *S$_i$* do
05.          use *t* to probe *R$_i$* for matches

# Oracle PGA Memory Management
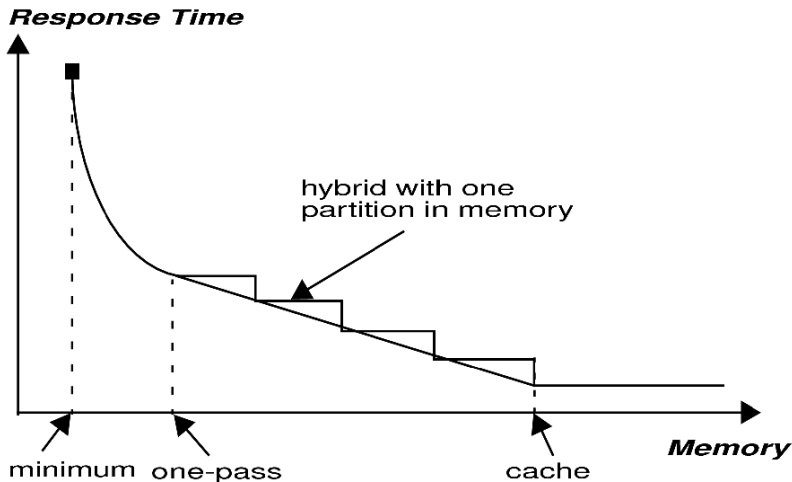
- **operator work area** = memory allocated for operator's execution
- Oracle's classification of work area size:
  - cache size = minimum work area size that enables evaluation to incur a single scan of input data
    - ⋆ cache size = size of input data + auxiliary memory structures
  - one-pass size = minimum work area size that enables evaluation to incur 2 scans of input data
  - multiple-pass size = work area size that is less than one-pass size
- cache size > one-pass size > multiple-pass size

# Effect of Memory on Sort Performance



**Response Time**

cache

**Memory**

minimum(64KB)    one-pass

(Dageville & Zait, VLDB 2002)

# Effect of Memory on Hash Join Perf.



(Dageville & Zait, VLDB 2002)

# Memory Management Approaches

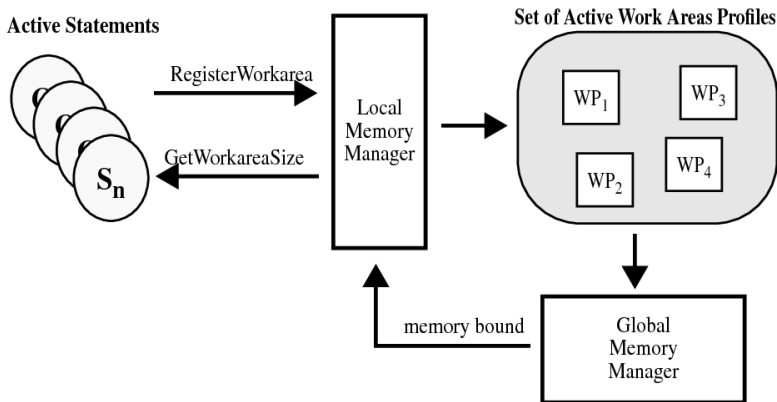How much memory to allocate to SQL operators?

- ► Fixed amount of memory for each operator
- ► Amount based on size of input operands
- ► Amount based on size of input operands + current workload
- ► Memory usage adapts to memory demand of system

# DBMS Approaches (circa 2002)

| DBMS | Initial work area size | Size during execution | Response to memory pressure |
|------|------------------------|-----------------------|------------------------------|
| Oracle 8i | static | static | none |
| **Oracle 9i** | **dynamic** | **adaptive** | **adaptive** |
| DB2/UDB 7.1 | static | static | minimum |
| Informix 9.3 | static | static | limit ops |
| SQL Server 7 | dynamic | static | queueing |
| Teradata | dynamic | static | ? |

# Memory Management Feedback Loop



(Dageville & Zait, VLDB 2002)

# Work Area Profile

Work area profile = metadata for an operator work area:

- ► Operator type
  - ► Example: sort, hash-join
- ► Current memory requirement to run with minimum/one-pass/cache memory
- ► Number instances of work area
  - ► ie degree of parallelism of operator
- ► Current amount of PGA memory

# Memory Management Feedback Loop

**SQL Operator**

1. Registers work area profile with local memory manager

5. Adjusts memory usage based on expected work area size (derived by LMM) & updates work area profile

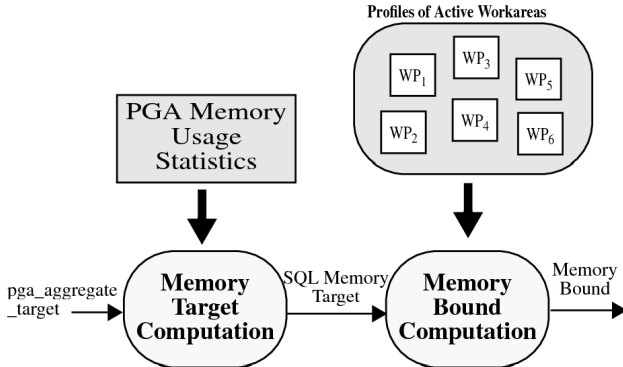**Local Memory Manager (LMM)**

2. Maintains active work area profiles in SGA

4. Computes expected work area size for each operator based on

- operator's work area profile
- memory bound (derived by GMM)

**Global Memory Manager (GMM)**

3. Computes memory bound for work area based on

- `pga_aggregated_target`
- PGA memory usage statistics
- active work area profiles

# GMM: Memory Bound Computation



Profiles of Active Workareas

(Dageville & Zait, VLDB 2002)

► **GMM** computes memory bound for work area based on

  ► `pga_aggregated_target`
  ► PGA memory usage statistics
  ► active work area profiles

► SQL memory target = amount of PGA allocated to work areas

# Memory Bound Computation

- Find largest value of $B$ such that

$$\sum_{i=1}^{N} ExpectedWorkAreaSize(WP_i, B) \leq SQL\ Memory\ Target$$

- $N$ = number of work area profiles
- $WP_i$ = memory profile of $i^{th}$ work area, $i \in [1, N]$
- Work area profile
  - minimum, one-pass & cache memory thresholds of work area

# Expected Work Area Size, $E$

- Computation of $E$ depends on type of operator (sort or non-sort) & execution mode (serial or parallel)

- $E \in$ [minimum threshold, cache threshold]

- **Type of operator**:

  - Operator is sort

    $$E = \begin{cases} \text{cache threshold} & \text{if } B \geq \text{ cache threshold,} \\ \text{one-pass threshold} & \text{if } B \geq \text{ one-pass threshold,} \\ \text{multi-pass threshold} & \text{otherwise.} \end{cases}$$

  - Operator is non-sort
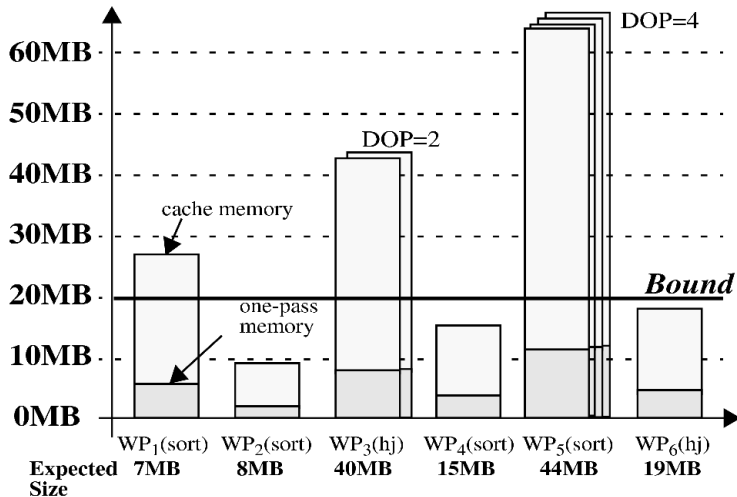
    $$E = \min\{B, \text{ cache threshold}\}$$

- **Execution mode**:

  $$E = \begin{cases} \min\{E, 0.05 \times \text{memory target}\} & \text{if serial,} \\ \min\{E \times DOP, 0.3 \times \text{memory target}\} & \text{otherwise.} \end{cases}$$

# Example of Memory Bound Computation

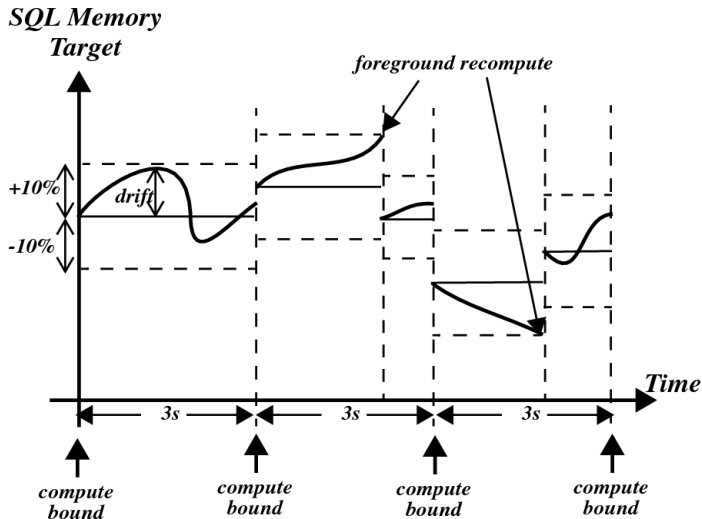(Dageville & Zait, VLDB 2002)



**Target = 133MB => Bound = 20M**

# Drift Management

- **GMM** is a background daemon
  - Recomputes memory bound periodically
- Limitation of background recomputation:
  - Slow to react to workload changes
- Memory drift
  - Measures staleness of memory bound
  - Expected amount of memory acquired/released by queries since last recomputation of memory bound
  - Drift can be positive or negative
- Triggers foreground recomputation when drift exceeds 10% of SQL memory target

# Foreground Computation of Bound



(Dageville & Zait, VLDB 2002)

# Memory Adaptive External Sorting

- **Creation of initial sorted runs**
  - Adapt size of sorted runs
- **Merging sorted runs**
  - Adapt merging fan-out
- Adapt size of input/output buffers

# Memory Adaptive Hash Join

- **Partitioning of build relation**
    - Adapt size of free buffer
    - Convert resident partitions to spilled ones to release memory
- **Partitioning of probe relation**
    - Adapt size of free buffer
    - Convert resident partitions to spilled ones to release memory
    - Restore spilled partitions to resident ones using extra memory
- Adapt size of input/output buffers

# References

**Required Readings**

- B. Dageville, M. Zait, SQL Memory Management in Oracle9i, VLDB 2002, 962-973.

**Additional Readings**

- W. Zhang, P.A. Larson, *Dynamic Memory Adjustment for External Mergesort*, VLDB 1997.

- H.H. Pang, M.J. Carey, M. Livny, *Partially preemptible hash joins*, SIGMOD 1993.

- A.J. Storm, C. Garcia-Arellano, S.S. Lightstone, Y. Diao, M. Surendra, *Adaptive self-tuning memory in DB2*, VLDB 2006.