# Class 16: Joins II

## Last time

### Nested-Loop Joins

| | | |
|---|---|---|
| Simple | $(P_R \cdot M) \cdot N + M$ | w/ R outer |
| Page-oriented | $M \cdot N + M$ | |
| Block-based | $\dfrac{M \cdot N}{K} + M$ | w/ K buffer |
| Index | $M + M \cdot P_R \cdot (\text{index\_access\_cost} + \text{data\_access\_cost})$ | |

hash ↙ ↘ $B^+$-tree        ↓clustered        unclustered
~1.2      2-4               1 I/o per page of
                            matching tuples

                                              1 I/o per
                                              matching tuple

## Sort-Merge Joins

$3 \cdot (M+N)$   if   $B > \sqrt{M}$   where M is # pages of the larger relation

$M+N$   if   $B > N$   where N corresponds to the smallest relation

## Today

→ Hash Joins
→ General Join Conditions
→ Aggregates

## Hash Joins

→ Use a hash function __h__ to create partitions of both relations   [hashing (building)]

→ match tuples only between the corresponding partitions
   [Probing (matching)]

B buffers  $R \bowtie S$
h hash function   $i=j$

**building** {
$\forall r \in R$
  read r and add it to buffer $h(r_i)$

$\forall s \in R$
  read s and add it to buffer $h(s_j)$
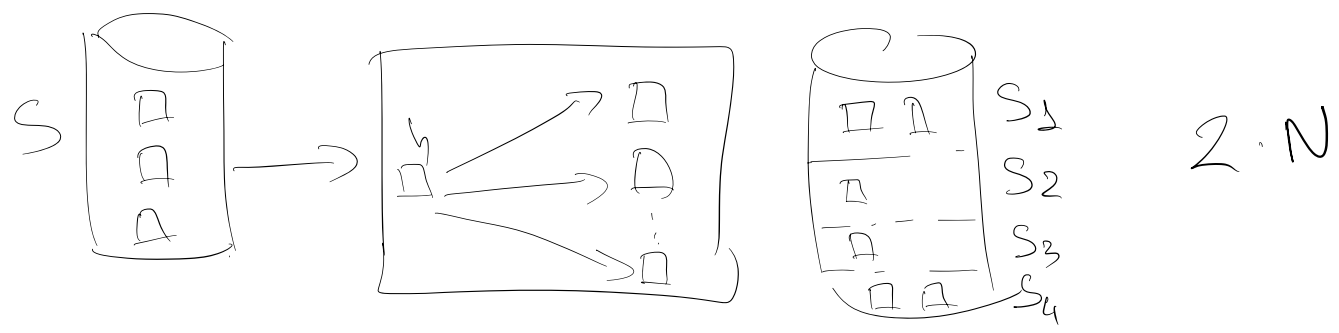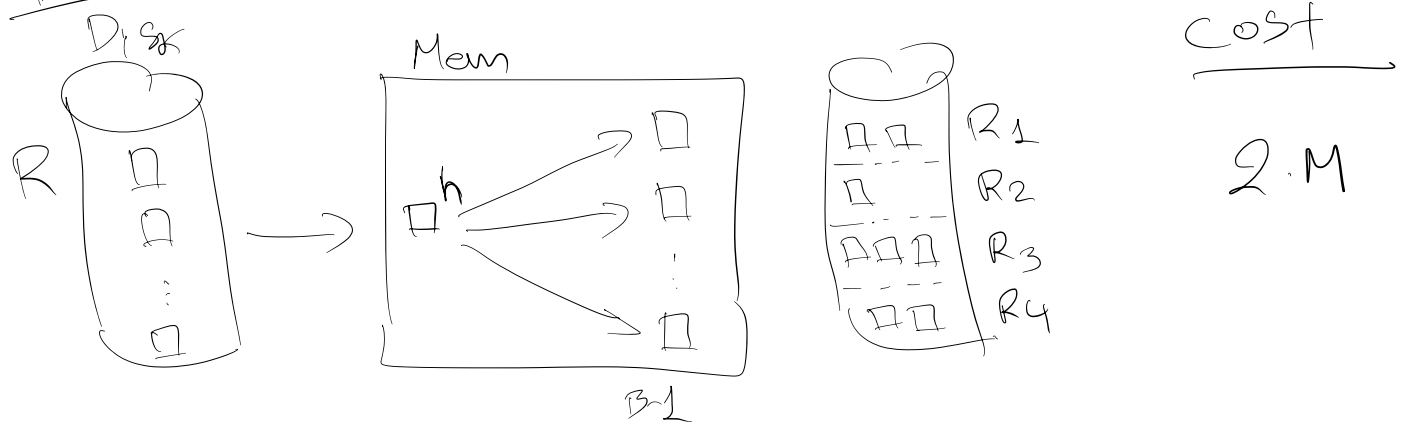}

**matching** {
for $l = 1, 2, \ldots K$

  $\forall r \in R_l$
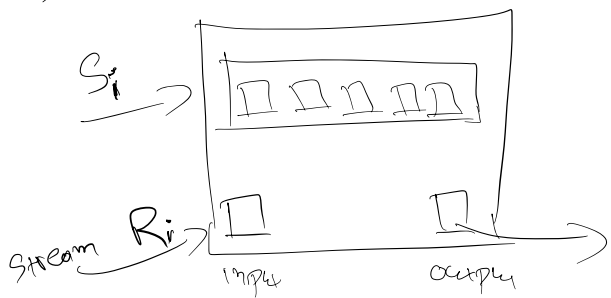    read r and insert into in-memory HT using $h_2(r_i)$

  $\forall s \in S_l$
    read s and probe HT using $h_2(s_j)$
    if match found add $\langle r, s \rangle$ to the result

  clear hash table from memory to proceed
  with next pair of partitions
}

Building


Disk   Mem   R1 R2 R3 R4   B-1

S   S1 S2 S3 S4

Cost
2·M

2·N

## Matching

S_i → [□ □ □ □ □]

Stream R_i → [□ (input)  □ (output)] →

read every partition once

in-memory HT w/ h2 (≠h)

Search in S_i as we stream R_i

Cost: $M + N$

total cost of Hash Join = $3(M+N) = \boxed{4500} → \boxed{4S}$

## Memory Requirements

→ enough buffer for the largest partition of the smaller relation (S)

→ Input page for the other relation

→ Output page

→ a few pages of hash metadata

Fudge factor $f$ (for example $f=1.04$)

if $h →$ uniform

size of a partition $\sim \dfrac{N}{B-1}$

$$B > \frac{f \cdot N}{B-1} + 2 \approx\!\Rightarrow B > \sqrt{f \cdot N}$$

what if not enough memory? (for S_i to fit in memory)
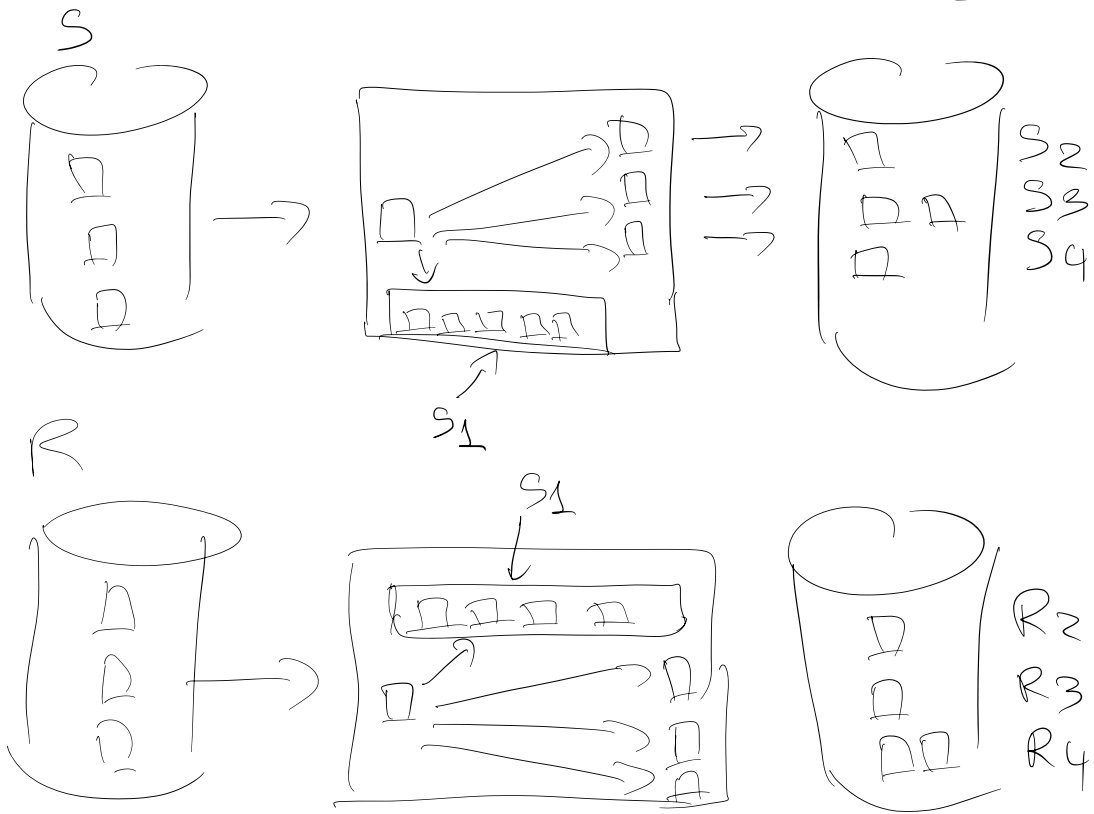
→ apply the same algorithm recursively

→ read, repartition S_i, R_i with h3 (≠h2, ≠h)

→ matching per subpartition (mem. is enough)

＊if not, again recursion

# What if we have more memory?

## Hybrid Hash Join



S → hashing → buckets → $S_2$, $S_3$, $S_4$ (with $S_1$ in memory)

R → hashing → $S_1$ in memory → probe → $R_2$, $R_3$, $R_4$

## Cost

| | | |
|---|---|---|
| → hashing | $S$ | $N + N - \text{sizeof}(S_1)$ |
| → hashing | $R$ | $M + M - \text{sizeof}(R_1)$ |
| → matching | | $M - \text{sizeof}(R_1) + N + \text{sizeof}(S_1)$ |

$$\text{total} \qquad 3(M+N) - 2\left(\text{sizeof}(S_1) + \text{sizeof}(R_1)\right)$$

$B = 300$
$M = 1000$
$N = 500$

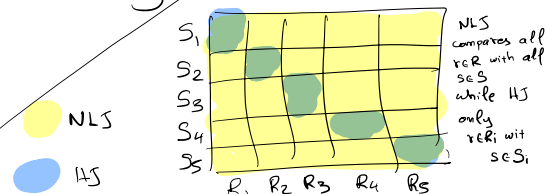$$3(1000 + 500) - 2(500 + 250) = 4500 - 1500 = \boxed{3000}$$

$\boxed{6s}$

if $B = 600$

read $S$ __once__ + build hash table

scan $R$ __once__ probe $S$ on-the-fly

## HJ vs NLJ

Work done during __matching__



NLJ compares all $r \in R$ with all $s \in S$ while HJ only $r \in R_i$ wit $s \in S_i$

NLJ

HJ

# Hash Join vs SMJ

|  | Hash Join | SMJ |
|---|---|---|
| cost | $3(M+N)$ | $3(M+N)$ |
| memory requirement (m.r.) | $B > \sqrt{f \cdot N}$ ← smaller | $B > \sqrt{M}$ ← larger |
| example: | $B > \sqrt{1.04 \cdot 500} = 23$ | $B > \sqrt{1000} = 32$ |

(cost is $>3(M+N)$)

| | | |
|---|---|---|
| $\sqrt{f N} \leq B < \sqrt{M}$ | $\left\{\begin{array}{l}\text{HJ exploits additional mem}\\ 3(M+N) - 2(\text{sizeof}(R_1) - \text{sizeof}(S_1))\end{array}\right\}$ | needs additional passes |
| $\sqrt{M} \leq B < N$ | | $3(M+N)$ |
| $B > N$ | $M+N$ | $M+N$ |
|  |  | sorted |

output
(if input sorted    $3(M+N)$    $M+N$

BUT   sensitive to data skew

---

(a) equality joins on several attributes

(b) inequality joins

→ (a) for INLJ we need index with all
       attributes in join conditions
   → sort/hash use combination of all attributes
→ (b) INLJ w/ $B^+$-Tree (not Hash Index)
       HJ/SMJ cannot work
       Block NLJ the best approach

## Set

UNION /EXCEPT (set difference)
   → Sorting
       → sort S+R on all attributes
   → merging → discard duplicates (UNION)
             → set-difference
   → refinement also applies

→ hashing
  → partition R+S
  → ∀ S-part probe corr. R-part
    → discard duplicates (UNION)
    → set_difference

→ Intersection → Special case of Join
    Equality across all attributes

---

Aggregation
  → SELECT AVG(sal) FROM E
      → SCAN once

  → GROUP BY
      ⟨age, avg_salary⟩

    hash (age) ⟶ ⟨age, salary, count⟩

    sort (age) Calculate "running info" of aggregation
                                        on-the-fly

  → if we have an index on ⟨Group-by, select, where⟩
        can use only the index [WAY FASTER]

  → Buffering
        #many thing in parallel
          tough to estimate what is costed by BP
  SNLJ  B>N ✓
        B<N    LRU → sequential flooding.
              MRU ✓.
  BNLJ replacement policy has no impact
  INLJ → sort the outer relation