# Class 19: 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 | $M + M \cdot P_R \cdot (index\_access\_cost + data\_access\_cost)$ | |

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
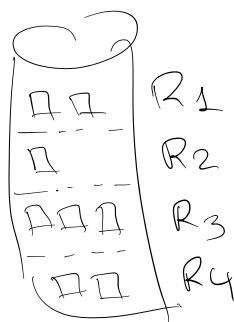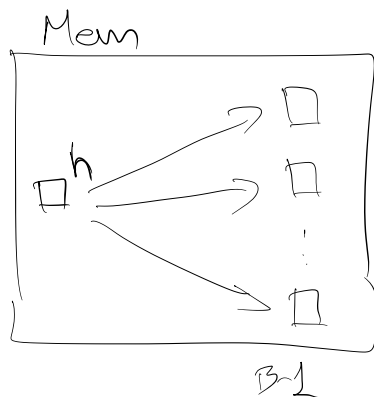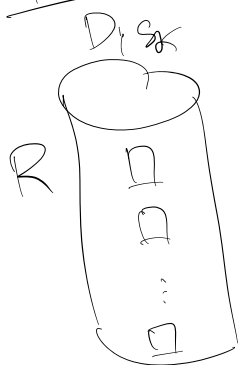[probing (matching)]

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

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

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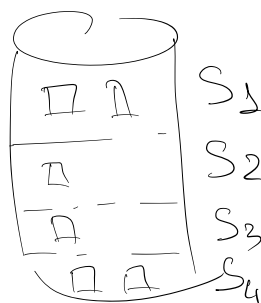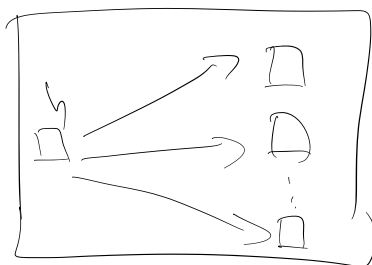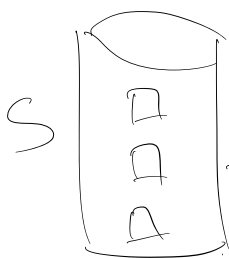read s and add it to buffer $h(s_j)$

matching
for $l = 1, 2, \dots K$

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

$$\forall s \in S_l$$
read s and probe AT 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

<u>Building</u>

Disk      Mem



R                    $R_1$
                     $R_2$
                     $R_3$
                     $R_4$
            B-1

<u>Cost</u>

$2 \cdot M$



S                    $S_1$
                     $S_2$
                     $S_3$
                     $S_4$
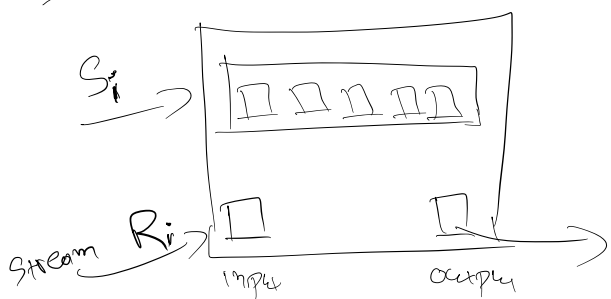
$2 \cdot N$

# Matching



read every partition once

in-memory HT w/ h2 ($\neq$ h)

Search in $S_i$ as we stream $R_i$

Cost: $M+N$

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

# 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 $h \to$ uniform

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

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

What if not enough memory? (for $S_i$ to fit in memory)
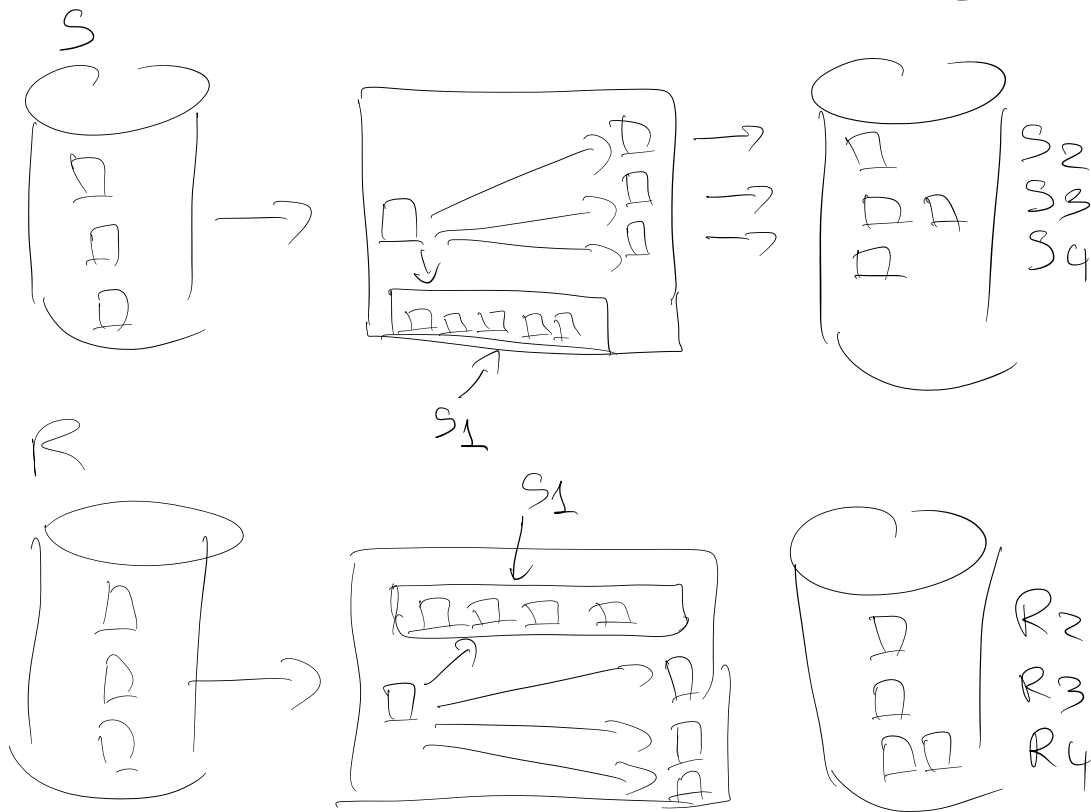
→ apply the same algorith recursively

→ read, repartition $S_i, R_i$ with $h_3$ ($\neq h_2, \neq h$)

→ matching per subpartition (mem. is enough)

✳ if not, again recursion

# What if we have more memory?

**Hybrid Hash Join**



## Cost

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

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

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

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

$\boxed{GS}$

if $B = 600$

read $S$ __once__ + build hash table

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

$$\text{Hash Join} \quad vs \quad SMJ$$

| | Hash Join | SMJ |
|---|---|---|
| cost | $3(M+N)$ | $3(M+N)$ |
| mem. req. | $BD \sqrt{f \cdot N} \leftarrow \text{smaller}$ | $BD \sqrt{M} \leftarrow \text{larger}$ |
| | $BD \sqrt{1.04 \cdot 500} = 23$ | $BD \sqrt{1000} = 32$ |
| $B(>m.r.) < N$ | $3(M+N) - 2(\text{sizeof}(R_1) - \text{sizeof}(S_1))$ | $3(M+N)$ |
| $B > N$ | $M+N$ | $M+N$ |
| | | sorted |
| output if input sorted | $\overline{3(M+N)}$ | $M+N$ |

$\underline{BUT}$     sensitive to data skew

---

(a) equality joins on several attributes

(b) inequality joins

$\rightarrow$ (a) for INLJ we need index with all attributes in join conditions

$\rightarrow$ sort/hash use combination of all attributes

$\rightarrow$ (b) INLJ w/ $B^+$-Tree (not Hash Index)

HJ/SMJ cannot work

Block NLJ the best approach

$\underline{Set}$

UNION/EXCEPT (set difference)

$\rightarrow$ sorting

$\rightarrow$ sort S+R on all attributes

$\rightarrow$ merging $\rightarrow$ discard duplicates (UNION)

$\rightarrow$ set-difference

→ 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 ✓
  INLJ → sort the outer relation