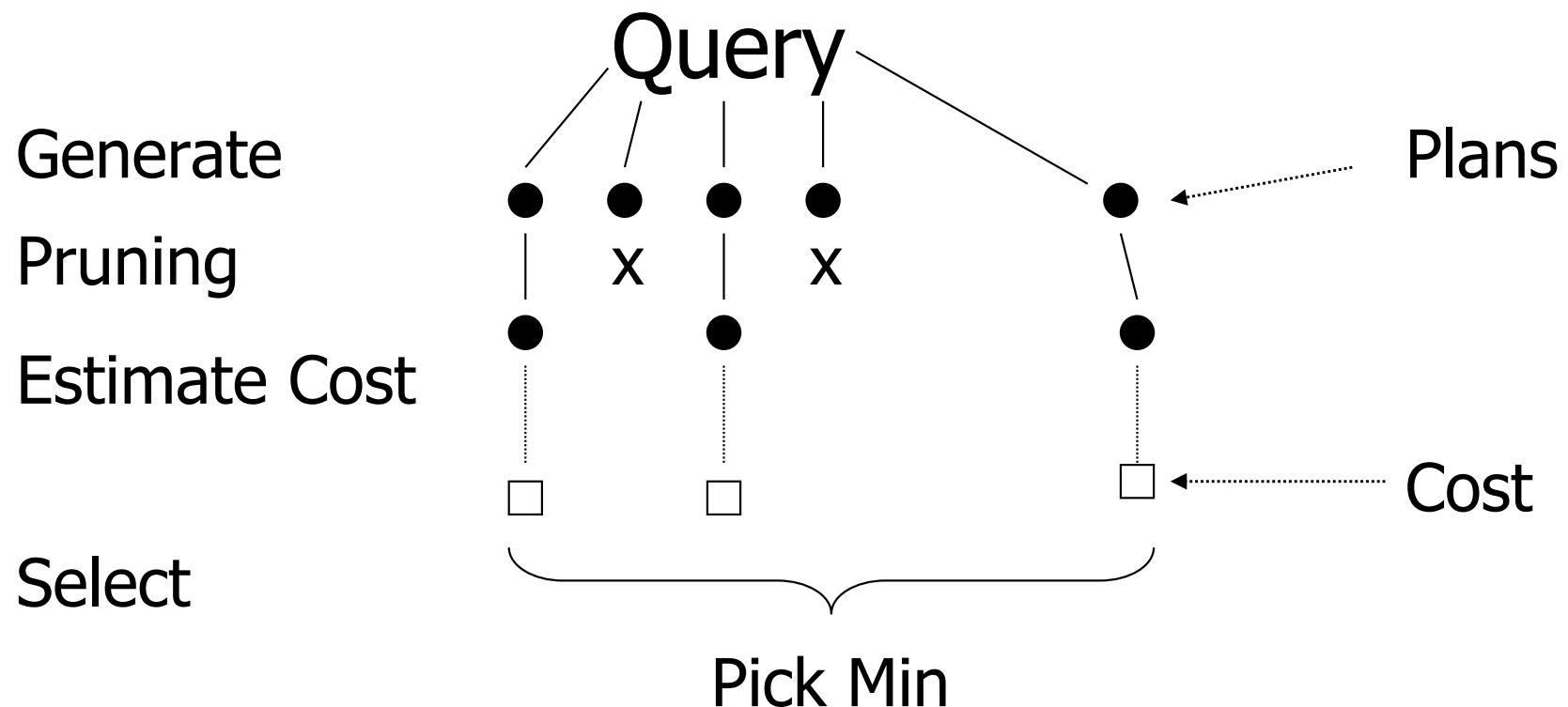


Query Optimization – Physical Query Plans

Hector Garcia-Molina and
Mahmoud Sakr

Query Optimization

--> Generating and comparing plans



To generate plans consider:

- Transforming relational algebra expression
(e.g. order of joins)
- Use of existing indexes
- Building indexes or sorting on the fly

- Implementation details:
 - e.g. - Join algorithm
 - Memory management
 - Parallel processing

Estimating IOs:

the disk access

by far the most restraining
part when doing
a query compared to CPU, -



- Count # of disk blocks that must be read (or written) to execute query plan

To estimate costs, we may have additional parameters:

$B(R)$ = # of blocks containing R tuples

$f(R)$ = max # of tuples of R per block

M = # memory blocks available

$T(R)$ = # of tuples in a relation (all tuples)

$T(R) > B(R)$

the number of tuples themselves

Notes 7

if store contiguously

the number of blocks containing the tuples

To estimate costs, we may have additional parameters:

$B(R)$ = # of blocks containing R tuples

$f(R)$ = max # of tuples of R per block

M = # memory blocks available

→ Height (How many levels in the tree)

$HT(i)$ = # levels in index i

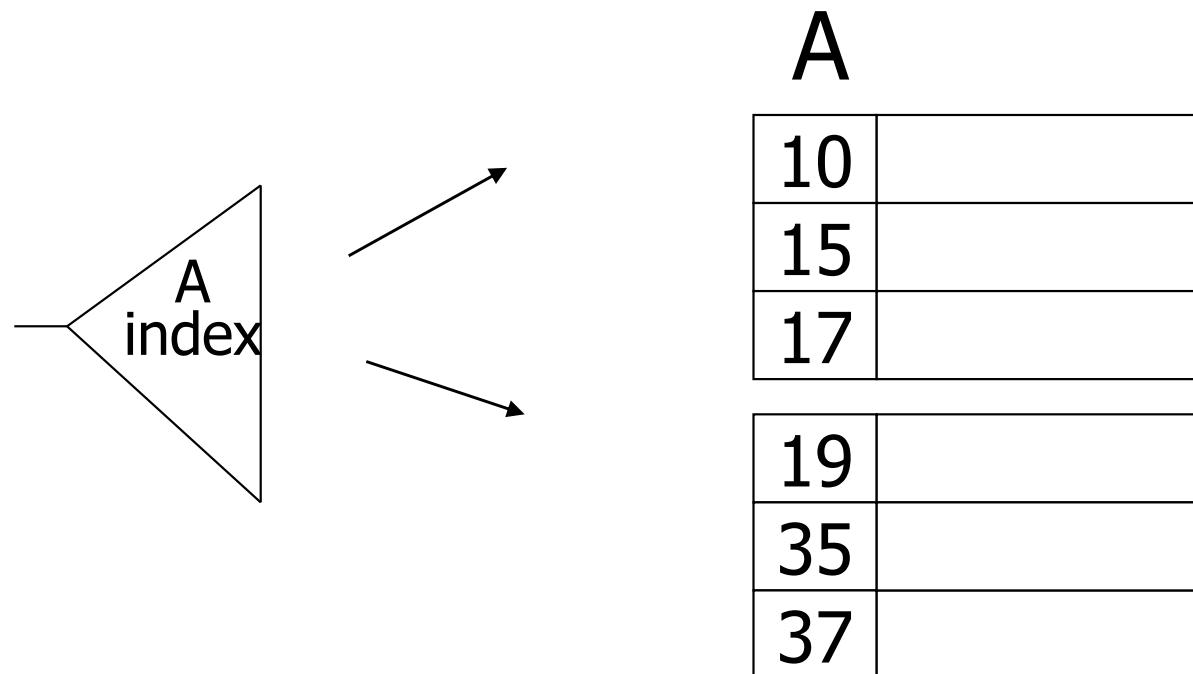
$LB(i)$ = # of leaf blocks in index i

using B-Tree indexes

Clustering index

organizing in order of
the key

Index that allows tuples to be read in an order that corresponds to physical order



Notions of clustering

blocks in the file can store tuples from multiple relations but they are ordered

- Clustered file organization

R1 R2 S1 S2

R3 R4 S3 S4

....

- Clustered relation $\Rightarrow B(R) \ll T(R)$

R1 R2 R3 R4

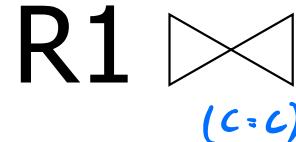
R5 R5 R7 R8

....

- Clustering index

Example $R1 \times R2$ over common attribute C

→ equijoin



$$T(R1) = 10,000$$

$$T(R2) = 5,000$$

$$S(R1) = S(R2) = 1/10 \text{ block}$$

$$\text{Memory available} = 101 \text{ blocks}$$

size of one tuple is 1 byte

1 block can store 10 tuples

⇒ we can store 1010 tuples
if contiguous in all the blocks
available

Example $R1 \bowtie R2$ over common attribute C

$T(R1) = 10,000$

$T(R2) = 5,000$

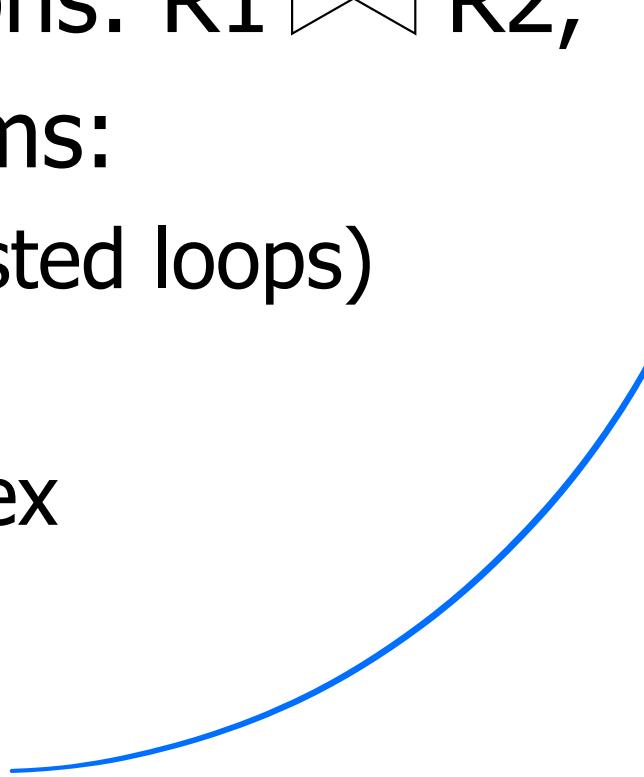
$S(R1) = S(R2) = 1/10$ block

Memory available = 101 blocks

→ Metric: # of IOs

(ignoring writing of result)

Options

- Transformations: $R1 \bowtie R2$, $R2 \bowtie R1$
 - Joint algorithms:
 - Iteration (nested loops)
 - Merge join
 - Join with index
 - Hash join
- 
- we can also invert
the order of the join*
- set of options
we can choose from*

- Iteration join (conceptually)

for each $r \in R_1$ do

two nested loops

for each $s \in R_2$ do

if $r.C = s.C$ then output r,s pair

and for each
possibility we output the
pair if it verifies
the join condition

- Merge join (conceptually) → by C value

(1) if R_1 and R_2 not sorted, sort them

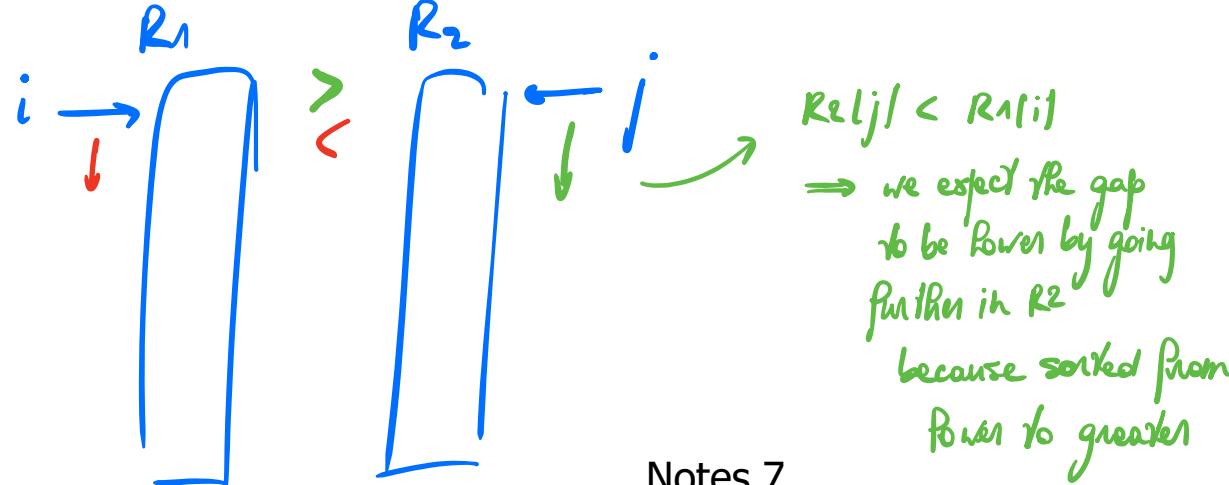
(2) $i \leftarrow 1; j \leftarrow 1;$

While ($i \leq T(R_1)$) \wedge ($j \leq T(R_2)$) do

 if $R_1\{ i \}.C = R_2\{ j \}.C$ then outputTuples

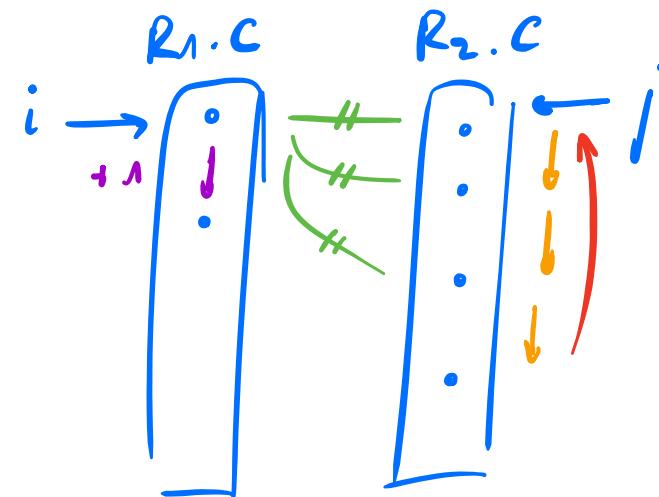
 else if $R_1\{ i \}.C > R_2\{ j \}.C$ then $j \leftarrow j+1$

 else if $R_1\{ i \}.C < R_2\{ j \}.C$ then $i \leftarrow i+1$



Procedure Output-Tuples

While $(R1\{ i \}.C \sqsubseteq R2\{ j \}.C) \wedge (i \leq T(R1))$ do
 $\underline{jj \leftarrow j;}$
 while $(R1\{ i \}.C = R2\{ jj \}.C) \wedge (jj \leq T(R2))$ do
 [output pair $R1\{ i \}, R2\{ jj \};$
 $jj \leftarrow \underline{\underline{jj+1}}$]
 $i \leftarrow \underline{\underline{i+1}}$]



Example

i	R1{i}.C	R2{j}.C	j
1	10	5	1
2	20	20	2
3	20	20	3
4	30	30	4
5	40	30	5
	→	50 ←	6
		52	7

- Join with index (Conceptually)

For each $r \in R1$ do

Assume R2.C index

[$X \leftarrow \text{index}(R2, C, r.C)$

for each $s \in X$ do

output r,s pair]

*we receive the tuples with
the same value as r.C or C*

Note: $X \leftarrow \text{index}(\text{rel}, \text{attr}, \text{value})$

then $X = \text{set of rel tuples with attr} = \text{value}$

- Hash join (conceptual)
 - Hash function h , range $0 \rightarrow k$
 - Buckets for R_1 : G_0, G_1, \dots, G_k
 - Buckets for R_2 : H_0, H_1, \dots, H_k

- Hash join (conceptual)
 - Hash function h , range $0 \rightarrow k$
 - Buckets for R_1 : G_0, G_1, \dots, G_k
 - Buckets for R_2 : H_0, H_1, \dots, H_k

Algorithm

- (1) Hash R_1 tuples into G buckets
- (2) Hash R_2 tuples into H buckets
- (3) For $i = 0$ to k do
 - match tuples in G_i, H_i buckets

Simple example

hash: even/odd

R1	R2
2	5
4	4
3	12
5	3
8	13
9	8
	11
	14

Even

2	4	8
R1		

4	12	8	14
R2			

Odd:

3	5	9
R1		

5	3	13	11
R2			

↳ Then we can join
buckets per buckets
⇒ smaller joins to compute

Factors that affect performance

- (1) Tuples of relation stored physically together?
- (2) Relations sorted by join attribute?
- (3) Indexes exist?

↗ nested Rof

Example 1(a) Iteration Join $R1 \bowtie R2$

(not sure if it is)

- Relations not contiguous \Rightarrow tuples are not in the same blocks
- Recall $\left\{ \begin{array}{l} T(R1) = 10,000 \quad T(R2) = 5,000 \\ S(R1) = S(R2) = 1/10 \text{ block} \\ \text{MEM}=101 \text{ blocks} \end{array} \right.$

Example 1(a) Iteration Join $R1 \bowtie R2$

- Relations not contiguous

- Recall $T(R1) = 10,000$

$$\left\{ \begin{array}{l} T(R1) = 10,000 \\ S(R1) = S(R2) = 1/10 \text{ block} \\ \text{MEM} = 101 \text{ blocks} \end{array} \right.$$

worst case: when reading next tuple we need to bring the next page into the memory

$$T(R2) = 5,000$$

Cost: for each $R1$ tuple:

[Read tuple + Read $R2$]

$$\text{Total} = 10,000 [1+5000] = 50,010,000 \text{ IOs}$$

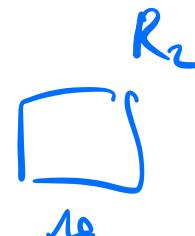
*For all tuples
of $R1$*

*read whole $R2$ must be
the tuple Notes 7 read to compare
 $R1$ to com with the $R1$ tuple*

- Can we do better?

Use our memory

- (1) Read 100 blocks of R1
- (2) Read all of R2 (using 1 block) + join
- (3) Repeat until done



$$(1000 + 5000) \times 10 = 60000$$

$> 50\,000\,000$

= 55000

Cost: for each R1 chunk:

Read chunk: 1000 IOs

Read R2: + $\frac{5000}{6000}$ IOs

what if we change the
memory allocated to each table:

$\frac{50}{R1} + \frac{51}{R2}$

$(\frac{500}{R1} + \frac{5000}{R2}) \times 20$

we still
have to
read the
whole relation

Mixing factor

$\frac{R2}{R1} \approx 1$

$(1000 + 10000) \times \boxed{5}$

= 55000

> 60000

less chunks
to read
⇒ better

Cost: for each R1 chunk:

Read chunk: 1000 IOs

Read R2: $\frac{5000}{6000}$ IOs

$$\text{Total} = \frac{10,000}{1,000} \times 6000 = 60,000 \text{ IOs}$$

number of
chunks of
R1

read chunk

+ read whole R2

>< 50.000.000

= way better

- Can we do better?

- Can we do better?
 - ➔ Reverse join order: $R2 \bowtie R1$

$$\text{Total} = \frac{5000}{1000} \times (1000 + 10,000) =$$

$$5 \times 11,000 = 55,000 \text{ IOs}$$

Example 1(b) Iteration Join R2 \bowtie R1

- Relations contiguous

nb! of chunks:

Size of table
Memory allocated

constant in
regard to previous
example (not contiguous)

- $\frac{10000}{1000} \times (1000 + 500) = 60000$
- $\left\lceil \frac{10000}{1000} \right\rceil \times (100 + 50) = 600$

↳ blocks in disk are contiguous $\rightarrow 1 \text{ block} = 10 \text{ relation tuples}$

PIP in the memory only requires 100 I/O to retrieve 100 blocks $\Rightarrow 1000 \text{ tuples}$

becomes

$1 \text{ block} = 1 \text{ tuple}$
 \forall previously

$\Rightarrow 1 \text{ I/O return}$
 $10 \text{ tuples} \rightarrow 10 \text{ times less I/O}$

Example 1(b) Iteration Join R2 \bowtie R1

- Relations contiguous

Cost

For each R2 chunk:

Read chunk: 100 IOs

Read R1: $\frac{1000}{1,100}$ IOs

Total = 5 chunks \times 1,100 = 5,500 IOs

of 1000 tuples
100 contiguous blocks

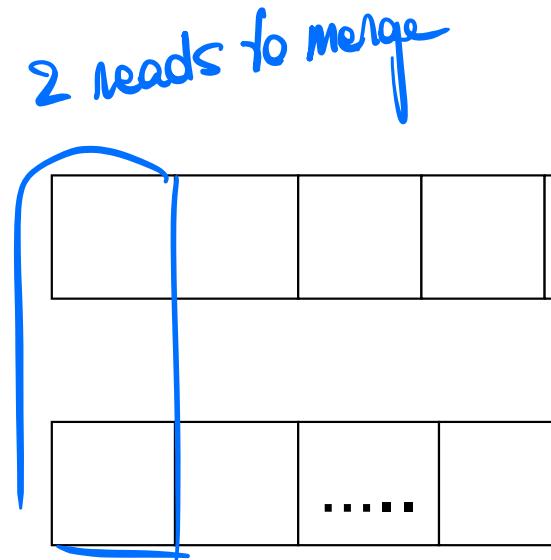
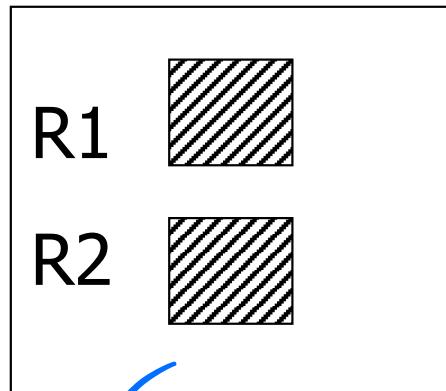
read
100 blocks of R2
+ read 1000 blocks of R1
↳ 10000 tuples

Example 1(c) Merge Join

→ attribute of equijoin

- Both R1, R2 ordered by C; relations contiguous

Memory



1000 tuples in
1000 pages
since contiguous

R1

R2

500 tuples in
500 pages

merge blocks, then advance from
one block either from R1 or R2
depending on which has lowest
value

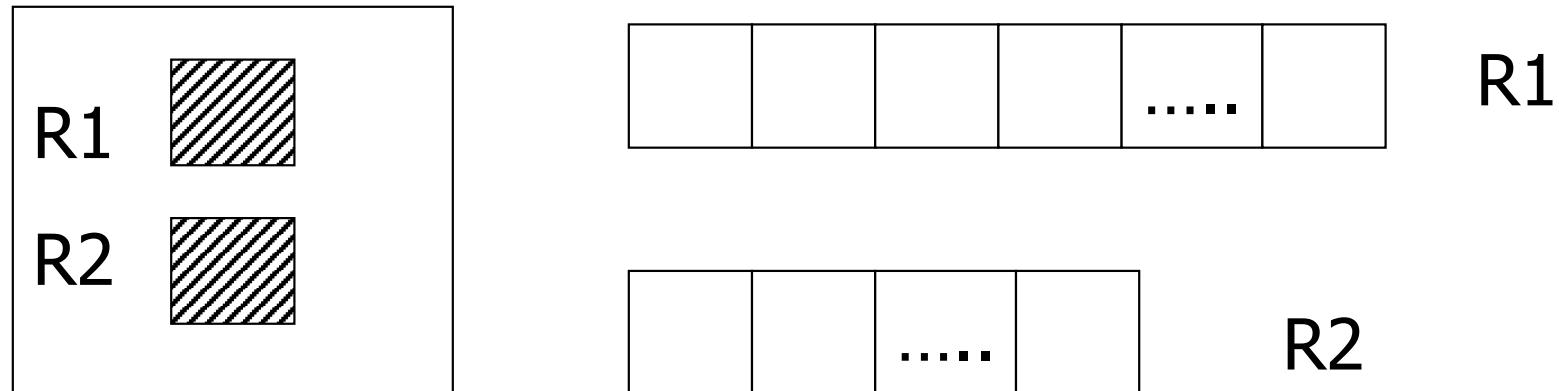
⇒ every block is read once

$$\Rightarrow 500 + 1000 = 1500 \text{ I/O}$$

Example 1(c) Merge Join

- Both R1, R2 ordered by C; relations contiguous

Memory



Total cost: Read R1 cost + read R2 cost
= 1000 + 500 = 1,500 IOs

Example 1(d) Merge Join

- R1, R2 not ordered, but contiguous

must be ordered for merge join

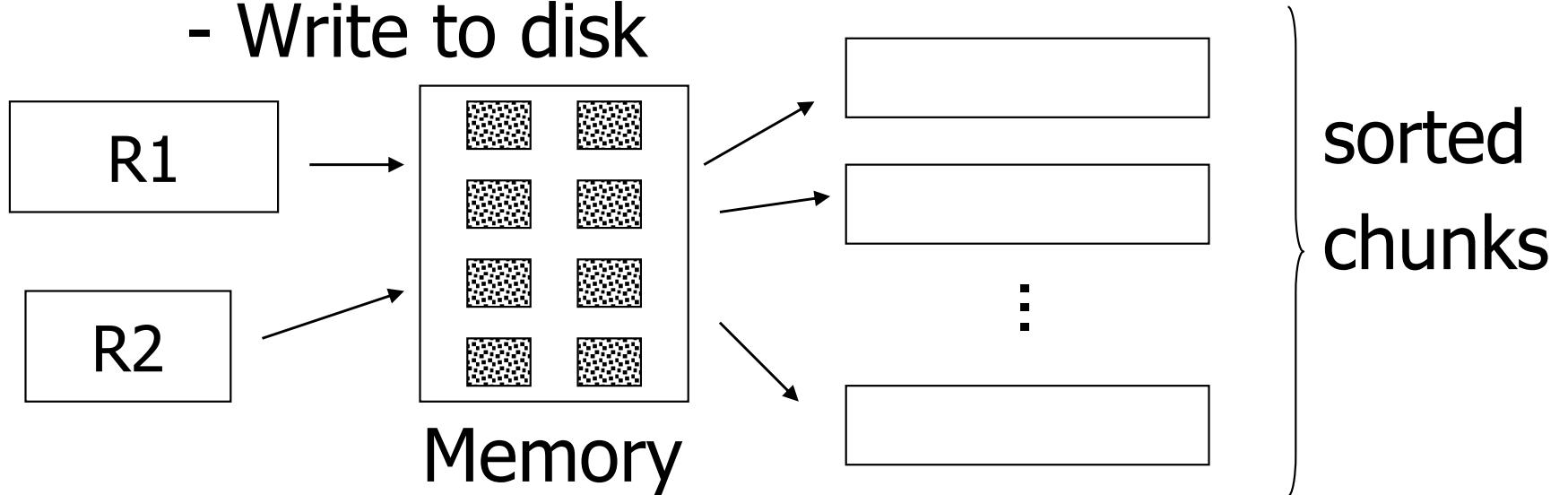
--> Need to sort R1, R2 first.... HOW?

One way to sort: Merge Sort

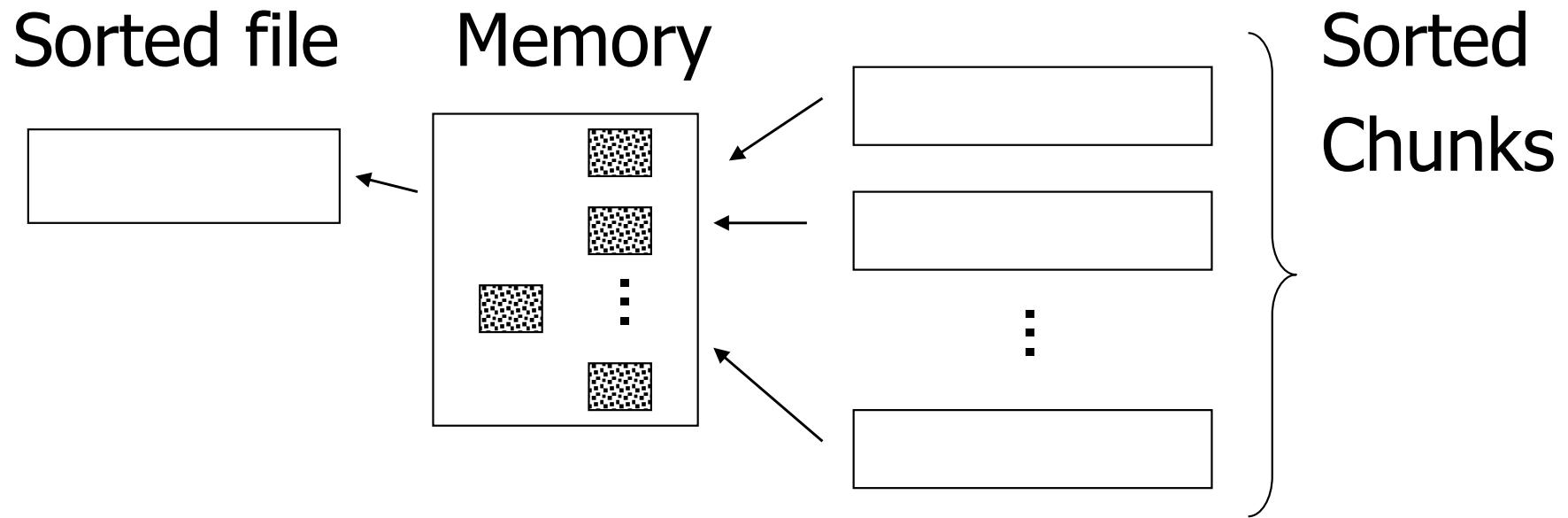
(i) For each 100 blk chunk of R:

- Read chunk
- Sort in memory
- Write to disk

*we sort each ~~repartio~~
independtly from the other*



(ii) Read all chunks + merge + write out



Cost: Sort

Each tuple is read,written,
read, written

so...

Sort cost R1: $4 \times 1,000 = 4,000$

Sort cost R2: $4 \times 500 = 2,000$

Merge Join vs Iteration Join on contiguous

a. Not contiguous

- 1. Iteration Join:
 - $R_1 \times R_2 : \frac{10.000}{1.000} \times (1000 + 500) = 60.000$
 - $R_2 \times R_1 : \frac{5.000}{1.000} \times (1000 + 10000) = 55.000$

2. Merge Join: require sequential + ordered

- Sequentiality:
 - $\underbrace{\text{read } R_1 + \text{write } R_1}_{\text{not seq.}} = 10000 + 1000 = 11.000$
 - $\underbrace{\text{read } R_2 + \text{write } R_2}_{\text{seq.}} = 5000 + 500 = 5.500$
- Order:
 - $2 \times (\text{read } R_1 + \text{write } R_1) = 2 \times 2000 = 4000$
 - $2 \times (\text{read } R_2 + \text{write } R_2) = 2 \times 1000 = 2000$
- join: $\text{read } R_1 + \text{read } R_2 = 1000 + 500 = 1.500 \Rightarrow \text{total: } 23.500$

b. Contiguous

- 1. Iteration Join:
 - $R_1 \times R_2 : \frac{10.000}{1.000} \times (100 + 50) = 6.000$

- $R_2 \times R_1 : \frac{5.000}{1.000} \times (100 + 100) = 5.500$

- 2. Merge Join:
 - order:
 - $2 \left(\underbrace{\text{read } R_1}_{1000} + \underbrace{\text{write } R_1}_{1000} \right) = 4000$
 - $2 \left(\underbrace{\text{read } R_2}_{500} + \underbrace{\text{write } R_2}_{500} \right) = 2000$

- join: $\text{read } R_1 + \text{read } R_2 = 1500 \Rightarrow \text{total} = 7500$

- c. Contiguous + ordered:
 - Iteration: $5.500 \rightarrow \text{nothing change}$
 - Merge: $1.500 \rightarrow \text{no need to sort first}$

$$\begin{aligned} \bullet T(R_1) &= 10000 \quad | \quad T(R_2) = 5000 \\ | \quad S(R_1) &= S(R_2) = 1/10 \\ | \quad \text{MEM} &= 101 \text{ blocks} \end{aligned}$$

Example 1(d) Merge Join (continued)

R1,R2 contiguous, but unordered

$$\begin{aligned}\text{Total cost} &= \text{sort cost} + \text{join cost} \\ &= 6,000 + 1,500 = 7,500 \text{ IOs}\end{aligned}$$

*Contiguous
and unordered*

$$\begin{aligned}\text{iteration : } & \frac{B(R_1)}{M-1} (M-1 + B(R_2)) \\ \text{merge : } & \left. \begin{array}{l} \cdot \text{sort } R_1 = 4B(R_1) \\ \cdot \text{sort } R_2 = 4B(R_2) \\ \cdot \text{merge} = B(R_1) + B(R_2) \end{array} \right\} = 5(B(R_1) + B(R_2))\end{aligned}$$

$$\frac{B(R_1)}{M-1} (M-1 + B(R_2)) = SB(R_1) + SB(R_2)$$

$$\Leftrightarrow B(R_1) \left(\frac{1}{M-1} (M-1 + B(R_2)) - S \right) = SB(R_2)$$

$$\Leftrightarrow B(R_1) \left(\frac{B(R_2)}{M-1} - 4 \right) = SB(R_2)$$

$$\Leftrightarrow B(R_1) = \frac{S B(R_2)}{\frac{B(R_2) - 4}{M-1}}$$

example : $M = 100$

$$B(R_2) = 5,000$$

$$\hookrightarrow \frac{25,000}{\frac{5000 - 4}{100}} = \frac{25,000}{46} \approx 544$$

iteration > merge

$$\text{if } B(R_1) < \frac{S B(R_2)}{\frac{B(R_2) - 4}{M-1}}$$

Example 1(d) Merge Join (continued)

R1,R2 contiguous, but unordered

$$\begin{aligned}\text{Total cost} &= \text{sort cost} + \text{join cost} \\ &= 6,000 + 1,500 = 7,500 \text{ IOs}\end{aligned}$$

But: Iteration cost = 5,500
so merge joint does not pay off!

But say $R_1 = 10,000$ blocks contiguous
 $R_2 = 5,000$ blocks not ordered

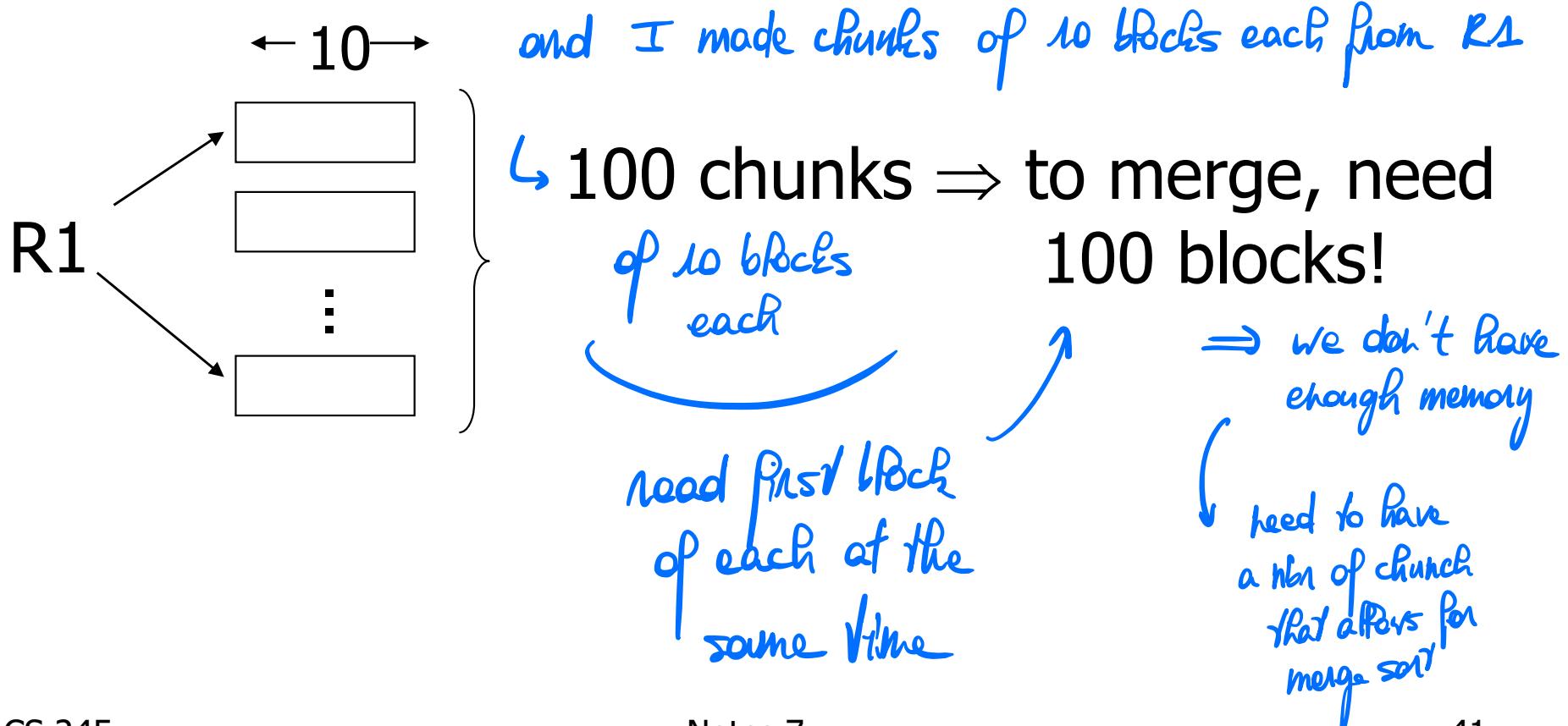
Iterate: $\frac{5000}{100} \times (100+10,000) = 50 \times 10,100$
 $= 505,000 \text{ IOs}$

Merge join: $5(10,000+5,000) = 75,000$ IOs

Merge Join (with sort) WINS!

How much memory do we need for merge sort?

E.g: Say I have 10 memory blocks (only)



In general:

Say k blocks in memory

x blocks for relation sort

chunks = (x/k) size of chunk = k

In general:

Say k blocks in memory

x blocks for relation sort

chunks = (x/k) size of chunk = k

chunks \leq buffers available for merge

In general:

Say k blocks in memory

x blocks for relation sort

chunks = (x/k) size of chunk = k

chunks \leq buffers available for merge

so... $(x/k) \leq k$

or $k^2 \geq x$ or $k \geq \sqrt{x}$

number of page
available
in memory

number of blocks
in the relation

In our example

R1 is 1000 blocks, $k \geq 31.62$ 32 needed

R2 is 500 blocks, $k \geq 22.36$ 23 needed

Need at least 32 buffers

Memory requirements of merge sort:

$$\cdot R_1 : T(R_1) = 10,000 \text{ tuples}$$

$$S(R_1) = 1/10 \text{ block/tuple} \Rightarrow f(R) = 10 \text{ tuples/block}$$

require 1000 contiguous blocks
to store R_1 entirely
 $= x$

$$\hookrightarrow \text{we want: } \# \text{ chunks} \leq \# \text{ blocks in memory} \Leftrightarrow \frac{x}{\ell} \leq \ell$$

$$= \ell$$

$$\Leftrightarrow \ell^2 \geq x$$

$$\Leftrightarrow \ell \geq \sqrt{x}$$

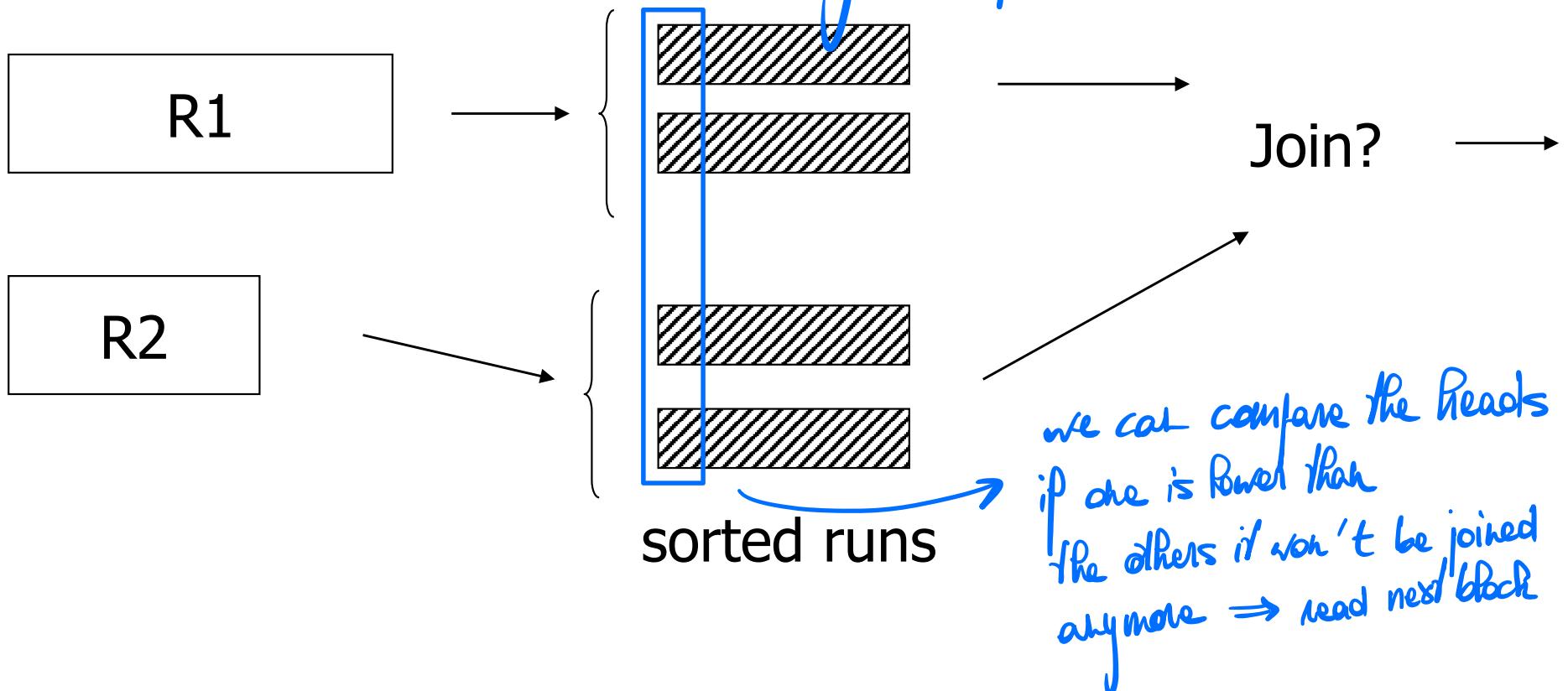
$$\hookrightarrow \text{Here } \sqrt{8} \approx 32$$

$$\frac{\# \text{ blocks of } R}{\# \text{ blocks in memory}} = \frac{x}{\ell}$$

\Rightarrow require 32 blocks of mem.

Can we improve on merge join?

Hint: do we really need the fully sorted files? → *do we need to merge completely?*



Cost of improved merge join:

$$\begin{aligned} C &= \text{Read R1} + \text{write R1 into runs} \\ &\quad + \text{read R2} + \text{write R2 into runs} \\ &\quad + \text{join} \\ &= 2000 + 1000 + 1500 = 4500 \end{aligned}$$

we do not
merge them
anymore

--> Memory requirement?

we need one memory
buffer per sorted runs
to apply the merge join

Example 1(e) Index Join

- Assume R1.C index exists; 2 levels
- Assume R2 contiguous, unordered
- Assume R1.C index fits in memory

⇒ querying the index is free

↳ happy case

↳ read all tuples of R2
is the only cost

Cost: Reads: 500 IOs
for each R2 tuple:

*because 500 tuples
in blocks of 10 tuples
⇒ 50 blocks to read*

- probe index - free
- if match, read R1 tuple: 1 IO

 output of a join requires
to read a block of R1
⇒ 1 IO per match

What is expected # of matching tuples?

(a) say R1.C is key, R2.C is foreign key

then expect = 1 → maximum one matching value
↪ # of different values (distinct values) since keys are unique

(b) say $V(R1,C) = 5000$, $T(R1) = 10,000$

with uniform assumption

expect = $10,000/5,000 = 2$

What is expected # of matching tuples?

→ domain of C (possible range)

(c) Say $\text{DOM}(R_1, C) = 1,000,000$

$$T(R_1) = 10,000$$

with alternate assumption

$$\text{Expect} = \frac{10,000}{1,000,000} = \boxed{\frac{1}{100}} \quad \% \text{ of matching tuples}$$

Total cost with index join

(a) Total cost = $500+5000(1)1 = 5,500$

(b) Total cost = $500+5000(2)1 = 10,500$

(c) Total cost = $500+5000(1/100)1=550$



⇒ requires one more parameter
"the selectivity"

What if index does not fit in memory?

Example: say R1.C index is 201 blocks

- Keep root + 99 leaf nodes in memory
- Expected cost of each probe is

$$E = (0)\frac{99}{200} + (1)\frac{101}{200} \approx 0.5$$

cost if in memory *cost of bringing a page to memory (if not in memory)*

(0 because access in memory is free)

Notes 7

Total cost (including probes)

$$= 500 + 5000 \text{ [Probe + get records]}$$

$$= 500 + 5000 [0.5 + 2] \quad \text{uniform assumption}$$

$$= 500 + 12,500 = 13,000 \quad (\text{case b})$$

Total cost (including probes)

$$= 500 + 5000 \text{ [Probe + get records]}$$

$$= 500 + 5000 [0.5 + 2] \quad \text{uniform assumption}$$

$$= 500 + 12,500 = 13,000 \quad (\text{case b})$$

For case (c):

$$= 500 + 5000[0.5 \times 1 + (1/100) \times 1]$$

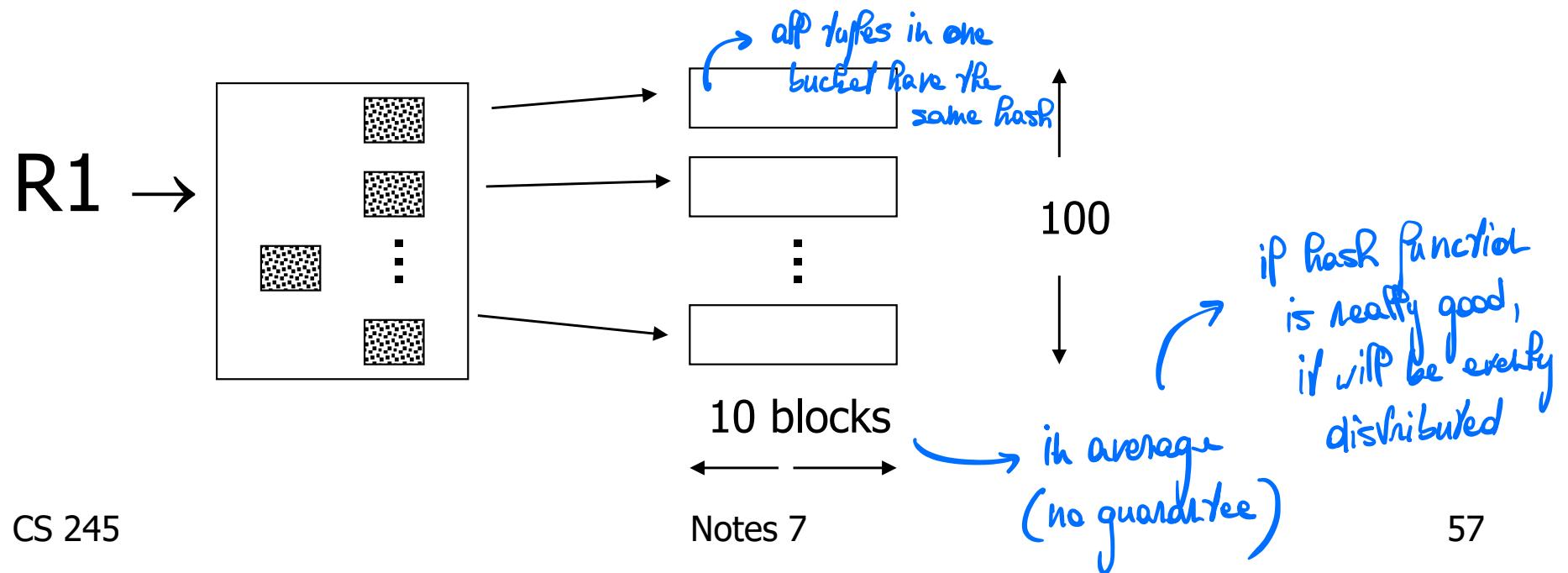
$$= 500 + 2500 + 50 = 3050 \text{ IOs}$$

So far

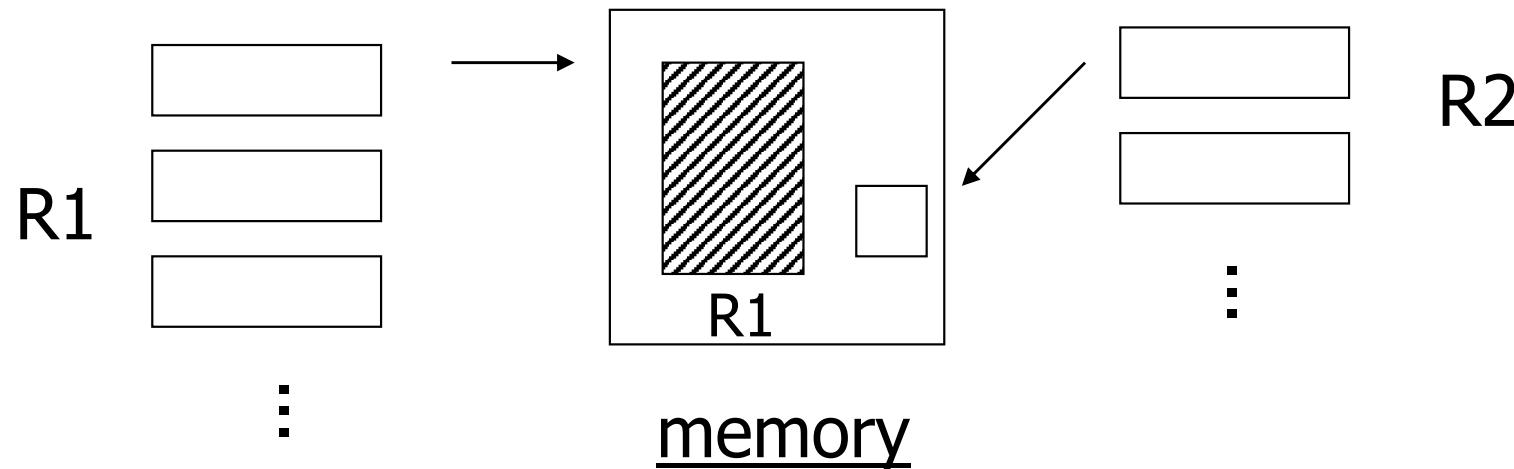
not contiguous	Iterate R2 \bowtie R1	55,000 (best)
	Merge Join	_____
	Sort+ Merge Join	_____
	R1.C Index	_____
	R2.C Index	_____
<hr/>		
contiguous	Iterate R2 \bowtie R1	5500
	Merge join	1500
	Sort+Merge Join	7500 \rightarrow 4500
	R1.C Index	5500 \rightarrow 3050 \rightarrow 550
	R2.C Index	_____

Example 1(f) Hash Join

- R1, R2 contiguous (un-ordered)
 - Use 100 buckets
 - Read R1, hash, + write buckets



- > Same for R2
- > Read one R1 bucket; build memory hash table
- > Read corresponding R2 bucket + hash probe



✉ Then repeat for all buckets

Cost:

“Bucketize:” Read R1 + write

 Read R2 + write

Join: Read R1, R2

Total cost = $3 \times [1000+500] = 4500$

Cost:

“Bucketize:” Read R1 + write

 Read R2 + write

Join: Read R1, R2

Total cost = $3 \times [1000+500] = 4500$

Note: this is an approximation since
buckets will vary in size and
we have to round up to blocks

Minimum memory requirements:

Size of R1 bucket = (x/k)

size of the buckets


k = number of memory buffers

x = number of R1 blocks

So... $(x/k) < k$

$k > \sqrt{x}$

need: $k+1$ total memory buffers

Memory Requirements:

- Iterative join :
 | 1 buffer for second relation read through
 | rest of them to hold part of the first relation

$$\rightarrow 101 \rightarrow 100/1$$
$$1001 \rightarrow 1000/1$$

↳ why so ? Let $T(R_1) = 10.000$ and $S(R_2) = S(R_1) = 1/10$

$$T(R_2) = 5.000$$

$$\hookrightarrow P(R_1) = P(R_2) = 10$$

→ with 101 MEM : $100/1$

↙

$$\frac{10.000}{100} \times (1000 + 5.000) = 6000 \times 10 = 60.000$$

Readings

- DATABASE SYSTEMS - The Complete Book, Hector Garcia-Molina, Jeffrey D. Ullman, Jennifer Widom, Second Edition.
- Chapter 16.7.1, 16.7.2, and relevant sections in Chapter 15.