## Discussion 6

**Iterators and Joins** 

#### Announcements

Vitamin 6 (Iterators and Joins) due Monday, March 4 at 11:59pm

Project 3 (Joins and QO) part 1 due Wednesday, March 6 at 11:59pm

#### Agenda

I. Iterators

#### II. Joins

- A. SNLJ
- B. PNLJ
- C. BNLJ
- D. INLJ
- E. Sort-Merge Join
- F. Grace Hash Join
- III. Worksheet

- Interface for **iterating** through data (making a single pass through the data)
- Important methods of an iterator are:
  - hasNext: is there another piece of data left
  - next: get the next piece of data
- Using an iterator is kind going through a linked list no support for random access
  - You can't (efficiently) say: fetch me the 50th item, then 30th, then 70th

- Recall: relational operators operate on relations and return relations
- We can implement this as: operate on an iterator (of the input relation) and return an iterator (of the output)
  - Optionally choose if we wish to materialize the output relation (write it to disk) or stream it to the next operator

- If we have:  $\pi_{id}(\sigma_{name > A'}(R))$ 
  - σ operator takes iterator over R, returns iterator that filters out tuples that don't satisfy predicate
  - $\circ~\pi$  operator takes in iterator from  $\sigma$  operator and returns tuples with only the id field
  - Materializing the relation returned from σ<sub>name > 'A'</sub>(R) not needed: π makes only one pass over the data
     Only need one page of R in memory at once

## Joins

#### Joins

- We'll be looking at inner (equi) joins
  - Algorithms can be pretty easily extended to left/right outer joins
  - Full joins require more thought not in scope
  - Some algorithms work for non-equi-joins, others don't
- A join is: taking one relation, and matching each tuple with tuples from another relation
- The join condition/predicate determines what rows in the other relation match to a row in the first relation

#### Joins

- Bit of notation:
  - $\circ$  [R] = number of pages in R
  - $\circ$  p<sub>R</sub> = number of records per page in R
  - |R| = number of records in R (the cardinality of R)
     |R| = p<sub>R</sub> \* [R]
- We typically *exclude* the final write's I/O cost
  - Don't add the cost of writing the joined output to disk
    - We might decide to stream it to the next operator instead of materializing results!

#### Simple Nested Loop Join (SNLJ)

- Direct translation of the definition of join into code
- To perform the join  $R \bowtie_{\theta} S$ , just take each row in R, and scan through S to find the matching rows!
  - for each row r in R:
    - for each row s in S:
      - if  $\theta(r, s)$ : output r joined with s



First iteration of outer loop...



Second iteration of outer loop...





Which relation should we pick as R and S respectively?

# $Cost = [R] + p_R[R][S]$

#### Worksheet Q1a

How many disk I/Os are needed to perform a simple nested loop join?

Companies: (company id, industry, ipo date) Nyse: (company\_id, date, trade, quantity)

- 20 pages of memory •
- Join Companies and NYSE on C.company id = N.company id
- company id is the primary key for Companies
- For every tuple in Companies, assume there are 4 matching tuples in NYSE
- [N] = 100 pages, p<sub>N</sub> = 100 tuples per page
  [C] = 50 pages, p<sub>C</sub> = 50 tuples per page
- Unclustered B+ indexes on C.company id and N.company id
- For both indexes, assume it takes 2 I/Os to access a leaf

#### Worksheet Q1a

How many disk I/Os are needed to perform a simple nested loop join?  $\underline{C \bowtie N}$ 

Cost is [C] + |C| \* [N] = [C] + p<sub>C</sub> [C] [N] = 50 + 50 \* 50 \* 100 = 250,050 I/Os

#### <u>N ⋈ C</u>

Cost is [N] + |N| \* [C] = [N] + p<sub>N</sub> [N] [C] = 100 + 100 \* 100 \* 50 = 500,100 I/Os Companies: (company\_id, industry, ipo\_date) Nyse: (company\_id, date, trade, quantity)

- 20 pages of memory
- We want to join Companies and NYSE on C.company\_id = N.company\_id
- company\_id is the primary key for Companies
- For every tuple in Companies, assume there are 4 matching tuples in NYSE
- [N] = 100 pages, p<sub>N</sub> = 100 tuples per page
- [C] = 50 pages,  $p_C = 50$  tuples per page
- Unclustered B+ indexes with height 1 on C.company\_id and N.company\_id

- Can we do better?
  - We scan S for every row in R, but we had to load an entire page of R into memory to get that row!
  - Instead of finding the rows in S that match a row in R, do the check for *all* rows in a page in R at once

- SNLJ
  - for each row r in R:
    - for each row s in S:
      - if  $\theta(r, s)$ : output r joined with s

- SNLJ (but with page fetches written out explicitly)
  - for each page  $P_R$  in R:
    - for each row r in  $P_R$ :
      - for each page P<sub>S</sub> in S:
        - for each row s in  $P_s$ :
          - if  $\theta(r, s)$ : output r joined with s

- PNLJ
  - for each page  $P_R$  in R:
    - for each page P<sub>s</sub> in S:
      - for each row r in P<sub>R</sub>:
        - for each row s in  $P_s$ :
          - if  $\theta(r, s)$ : output r joined with s



First iteration of outer loop...



First iteration of outer loop...



First iteration of outer loop...



Second iteration of outer loop...



Second iteration of outer loop...





#### Block Nested Loop Join (BNLJ)

- Can we do even better?
  - We only use three page of memory for PNLJ (one buffer for R, one buffer for S, one output buffer), but we usually have more memory!
  - Instead of fetching one page of R at a time, why not fetch as many pages of R as we can fit (B - 2 pages)!

#### Block Nested Loop Join (BNLJ)

- PNLJ
  - for each page  $P_R$  in R:
    - for each page P<sub>S</sub> in S:
      - for each row r in  $P_R$ :
        - for each row s in  $P_s$ :
          - if  $\theta(r, s)$ : output r joined with s

#### Block Nested Loop Join (BNLJ)

- BNLJ
  - for each block of B 2 pages  $C_R = \{P_1, P_2, \dots, P_{B-2}\}$ in R:
    - for each page P<sub>S</sub> in S:
      - for each row r in  $C_R$ :
        - for each row s in  $P_s$ :
          - if  $\theta(r, s)$ : output  $\vec{r}$  joined with s



	_
	_

R



R





R



B = 4

_		
_		



	_
R



#### Output Buffer

R



#### Output Buffer

_		_

R



B = 4

#### Output Buffer

		_

#### BNLJ

	page 0 record 0	page 0 record 1	page 0 record 2	page 0 record 3	page 1 record 0	page 1 record 1	page 1 record 2	page 1 record 3	page 2 record 0	page 2 record 1	page 2 record 2	page 2 record 3	page 3 record 0	page 3 record 1	page 3 record 2	page 3 record 3
page 0 record 0	_	-	-			-	-		_	_		_	_	-	I	-
page 0 record 1																
page 0 record 2						2-		_								
page 0 record 3																
page 1 record 0						22	8-3	2.3		1			-8			
page 1 record 1																
page 1 record 2																
page 1 record 3																
page 2 record 0			1			<u>9</u>	9 - 3 	1								1
page 2 record 1							F									
page 2 record 2						Č.										
page 2 record 3	$\vdash$		1		$\vdash$	1				-				-		
page 3 record 0	$\vdash$				$\vdash$	2	-			-						
page 3 record 1	$\vdash$															
page 3 record 2	⊢	-			$\vdash$	2	-	-		-		-	-	-		-
E broom E apen			1.1			1		1							1	

leftRecordIterator: Iterator over left block rightRecordIterator: Iterator over right page leftRecord: rightRecord :

### Cost of BNLJ?

[R] + (# blocks in R) \* [S] $= [R] + \Gamma[R] /_{chunksize} \Im * [S]$ 

=[R] + Γ[R]/<sub>(B - 2)</sub> <sup>¬</sup>[S]

#### Worksheet Q1b

How many disk I/Os are needed to perform a block nested loop join?

- 20 pages of memory
- We want to join Companies and NYSE on C.company\_id = N.company\_id
- company\_id is the primary key for Companies
- For every tuple in Companies, assume there are 4 matching tuples in NYSE
- [N] = 100 pages, p<sub>N</sub> = 100 tuples per page
- [C] = 50 pages,  $p_c = 50$  tuples per page
- Unclustered B+ indexes with height 1 on C.company\_id and N.company\_id

#### Worksheet Q1b

How many disk I/Os are needed to perform a block nested loop join?

B = 20, block size = B - 2 = 18

#### <u>C ⋈ N</u>

- Cost is [C] + [[C] / B 2] \* [N]
- = 50 + [50 / 18] \* 100 = 350 I/Os

#### <u>N ⋈ C</u>

Cost is [N] + [[N] / B - 2] \* [C]

= 100 + [100 / 18] \* 50 = 400 I/Os

Companies: (company\_id, industry, ipo\_date) Nyse: (company\_id, date, trade, quantity)

- 20 pages of memory
- We want to join Companies and NYSE on C.company\_id = N.company\_id
- company\_id is the primary key for Companies
- For every tuple in Companies, assume there are 4 matching tuples in NYSE
- [N] = 100 pages, p<sub>N</sub> = 100 tuples per page
- [C] = 50 pages,  $p_c = 50$  tuples per page
- Unclustered B+ indexes with height 1 on C.company\_id and N.company\_id

I/O cost for BNLJ: min(350, 400) I/Os = **350 I/Os** 

- A join is essentially:
  - for each row r in R:
    - for each row s in S that satisfies  $\theta(r, s)$ :
      - output r joined with s

- An *index on S* allows us to do the inner loop efficiently!
  - for each row r in R:
    - for each row s in S that satisfies θ(r, s)
       (found using the index):
      - output r joined with s

- What's the I/O cost?
  - [R] + |R| \* cost to find matching S tuples
    - [R] from scanning through R
  - Cost to find matching S tuples:
    - Alternative 1: cost to traverse root to leaf + read all the leaves with matching tuples
    - Alternative 2/3: cost of retrieving RIDs (similar to Alternative 1) + cost to fetch actual records
      - 1 I/O per page if clustered, 1 I/O per tuple if not

- What's the I/O cost?
  - $\circ$  [R] + |R| \* cost to find matching S tuples
    - [R] from scanning through R
  - If we have no index, then the only way to search for matching S tuples is by scanning all of S  $\rightarrow$  SNLJ, PNLJ, BNLJ, etc.
    - Cost to find matching S tuples is then [S], giving us the formula for SNLJ cost

R.col



R.col



R.col



R.col



R.col



#### Worksheet Q1c

How many disk I/Os are needed to perform an index nested loop join?

- 20 pages of memory
- We want to join Companies and NYSE on C.company\_id = N.company\_id
- company\_id is the primary key for Companies
- For every tuple in Companies, assume there are 4 matching tuples in NYSE
- [N] = 100 pages,  $p_N = 100$  tuples per page
- [C] = 50 pages,  $p_c = 50$  tuples per page
- Unclustered alternative 3 B+ indexes with height 1 on C.company\_id and N.company\_id. Throughout the problem assume no index nodes are cached.

#### Worksheet Q1c

How many disk I/Os are needed to perform an index nested loop join?

#### <u>C</u> ⋈ N

Cost is [C] + |C| \* cost of searching N = 50 + (50 \* 50) \* (2 + 4) = 15,050 I/Os

#### <u>N ⋈ C</u>

Cost is [N] + |N| \* cost of searching C = 100 + (100 \* 100) \* (2 + 1) = 30,100 I/Os

I/O cost: min(30,100, 15,050) = **15,050 I/Os** 

- 20 pages of memory
- We want to join Companies and NYSE on C.company\_id = N.company\_id
- company\_id is the primary key for Companies
- For every tuple in Companies, assume there are 4 matching tuples in NYSE
- [N] = 100 pages,  $p_N = 100$  tuples per page
- [C] = 50 pages, p<sub>c</sub> = 50 tuples per page
- Unclustered alternative 3 B+ indexes with height 1 on C.company\_id and N.company\_id. Throughout the problem assume no index nodes are cached.

#### Worksheet Q1d

Now assume the index on NYSE.company\_id is clustered. What is the cost of an index nested loop join using companies as the outer relation?

- 20 pages of memory
- We want to join Companies and NYSE on C.company\_id = N.company\_id
- company\_id is the primary key for Companies
- For every tuple in Companies, assume there are 4 matching tuples in NYSE
- [N] = 100 pages,  $p_N = 100$  tuples per page
- [C] = 50 pages,  $p_c = 50$  tuples per page
- Unclustered alternative 3 B+ indexes with height 1 on C.company\_id and N.company\_id. Throughout the problem assume no index nodes are cached.

#### Worksheet Q1d

Now assume the index on NYSE.company\_id is clustered. What is the cost of an index nested loop join using companies as the outer relation?

#### <u>C</u> ⋈ N

Cost is [C] + |C| \* cost of searching N

= 50 + 50 \* 50 \* (2 + # pages of matching tuples)

 $= 50 + 50 * 50 * (2 + ceil(# matches /p_N))$ 

= 50 + 50 \* 50 \* (2 + ceil(4/100)) = **7550 l/Os** 

- 20 pages of memory
- We want to join Companies and NYSE on C.company\_id = N.company\_id
- company\_id is the primary key for Companies
- For every tuple in Companies, assume there are 4 matching tuples in NYSE
- [N] = 100 pages,  $p_N = 100$  tuples per page
- [C] = 50 pages,  $p_C = 50$  tuples per page
- Unclustered alternative 3 B+ indexes with height 1 on C.company\_id and N.company\_id. Thoughout the problem assume no index nodes are cached.

### Sort-Merge Join (SMJ)

- What if we process the data a bit before we join things together?
  - For example, sort both relations first! Then we can join them efficiently
  - In some cases, we might even have one of the relations already sorted on the right key, and then we don't even have to spend time sorting it!

### Sort-Merge Join (SMJ)

- First step: **sort** both R and S (with external sorting)
- Second step: merge matching tuples from R and S together
  - We do this efficiently by moving iterators over sorted R and sorted S in lockstep: move the iterator with the smaller key
    - We know that this key is smaller than all remaining key values in the other relation, so we're completely done joining that tuple!

### Sort-Merge Join (SMJ)

- First step: **sort** both R and S (with external sorting)
- Second step: merge matching tuples from R and S together
  - Need a bit more care than this: we might have multiple rows in R matching with multiple rows in S
    - Mark the first matching row in S, match tuples with the first matching row in R, then reset the iterator to the mark so we can go through the rows in S again for the second matching row in R



<pre>while not done {</pre>	sid	sname	sid	bid
<pre>while (r &lt; s) { advance r }</pre>	22	dustin	━⇒28	103
<pre>while (r &gt; s) { advance s }</pre>	<b>28</b>	уирру	28	104
mark s // save start of "hlock"	31	lubber	31	101
while $(r == s)$ {	31	lubber2	31	102
// Outer loop over r	44	guppy	42	142
while $(r == s)$ {	57	rusty	58	107
// Inner loop over s				
yield <r, s=""></r,>				
advance s				
}				
<b>reset</b> s to mark				
advance r				
}				



sname	sid	bid
dustin	<b>→</b> 28	103
уирру	28	104
lubber	31	101
lubber2	31	102
guppy	42	142
rusty	58	107

bid

<pre>while not done {</pre>	sid	snar	ne	sid
while $(r < s)$ { advance r }	22	dust	in	<b>2</b> 8
while (r > s) { advance s }	<b>28</b>	yupp	у	28
mark < // save start of "block"	31	lubb	er	31
while $(r == s)$ {	31	lubb	er2	31
// Outer loop over r	44	gupp	ру	42
while $(r == s)$ {	57	rusty	/	58
// Inner loop over s				
yield <r, s=""></r,>		sid	sname	bid
advance <u>s</u>		28	yuppy	103
}				
<b>reset</b> s to mark				
advance r				
}				
}				



sid	snam	е	sid	bid
22	dustin	ì	28	103
28	yuppy	,	<b>→</b> 28	104
31	lubbe	r	31	101
31	lubbe	r2	31	102
44	guppy	,	42	142
57	rusty		58	107
	sid	sname	bid	
	28	yuppy	103	

bid

<pre>while not done {</pre>	si	id	snam	e	sid
while $(r < s) \{ advance r \}$	2	2	dusti	n	28
<pre>while (r &gt; s) { advance s }</pre>	<b>→</b> 2	8	yupp	у	₩28
mark < // save start of "block"	3	1	lubbe	er	31
while $(r == s)$ {	3	1	lubbe	er2	31
// Outer loop over r	4	4	gupp	у	42
while (r == s) {	5	7	rusty		58
// Inner loop over s					
yield <r, s=""></r,>		ę	sid	sname	bid
advance s			28	yuppy	103
}			28	yuppy	104
<b>reset</b> s to mark					
advance r					
}					
}					



	sid	sname	e	sid	bid
	22	dustin	I	28	103
_	28	yuppy		28	104
	31	lubbe	r	<b>→</b> 31	101
	31	lubbe	r2	31	102
	44	guppy	,	42	142
	57	rusty		58	107
		sid	sname	bid	
		28	yuppy	103	
		28	yuppy	104	

```
while not done {
while (r < s) { advance r }
while (r > s) { advance s }
mark s // save start of "block"
while (r == s) {
 // Outer loop over r
  while (r == s) {
   // Inner loop over s
   yield <r, s>
    advance s
  }
  reset s to mark
  advance r
```

sid	snan	ne	sid	bid
22	dusti	'n	≥28	103
28	yupp	y	28	104
31	lubb	er	31	101
31	lubb	er2	31	102
44	gupp	у	42	142
57	rusty	,	58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	

w <b>hile</b> not done {
while (r < s) { advance r }
while $(r > s)$ { advance s }
<pre>mark s // save start of "block"</pre>
while $(r == s)$ {
// Outer loop over r
while $(r == s)$ {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
<b>reset</b> s to mark
advance r
}

sid	sname		sid	bid	
22	dusti	า	<b>→</b> 28	103	
28	уирру		28	104	
31	lubbe	er	31	101	
31	lubbe	er2	31	102	
44	guppy	y	42	142	
57	rusty		58	107	
	sid	sname	bid		
	28	yuppy	103		
	28	yuppy	104		

bid

sid

<pre>while not done {</pre>	sid	snai	ne	sid
<pre>while (r &lt; s) { advance r }</pre>	22	dust	in	<b>→</b> 28
while (r > s) { advance s }	28	yup	ру	28
mark s // save start of "block" 📒	31	lubb	ber	31
while $(r == s)$ {	31	lubber2		31
// Outer loop over r	44	gup	ру	42
while $(r == s)$ {	57	rust	у	58
// Inner loop over s				
yield <r, s=""></r,>		sid	sname	bid
advance s		28	yuppy	103
}		28	yuppy	104
<b>reset</b> s to mark				
advance r				
}				

<pre>while not done {</pre>		snar	ne	sid	bid
while $(r < s)$ { advance r }	22	dust	in	28	103
while $(r > s)$ { advance s }	28	yupp	ру	<b>2</b> 8	104
mark s // save start of "block"	→31	lubber		31	101
while $(r == s)$ {	31	lubber2		31	102
// Outer loop over r	44	guppy rusty		42	142
while $(r == s)$ {	57			58	107
// Inner loop over s					
yield <r, s=""></r,>		sid sname		bid	
advance s		28	yuppy	103	
}		28	yuppy	104	
<b>reset</b> s to mark					
advance r					
}					

<pre>while not done {</pre>	sid	sna	me	sid	bid
while (r < s) { advance r }	22	dus	tin	28	103
<pre>while (r &gt; s) { advance s }</pre>	28	yup	ру	28	104
mank of (/ cave stant of "block"	<b>→</b> 31	lubk	ber	<b>→</b> 31	101
while $(r == s)$ {	31	lubb	per2	31	102
// Outer loop over r	44	gup	ру	42	142
while $(r == s)$ {	57	rust	.y	58	107
// Inner loop over s					
yield <r, s=""></r,>		sid	sname	bid	
advance s		28	yuppy	103	
}		28	yuppy	104	
<b>reset</b> s to mark					
advance r					
}					


bid

sid

<b>while</b> not done {	sid	snam	e	sid
while $(r < s) \{ advance r \}$	22	dusti	n	28
while (r > s) { advance s }	28	yupp	y	28
mark < // save start of "block"	31	lubbe	er	<b>&gt;</b> 31
while $(r == s)$ {	31	lubbe	er2	31
// Outer loop over r	44	gupp	y	42
while (r == s) {	57	rusty		58
// Inner loop over s				
yield <r, s=""></r,>		sid	sname	bid
advance <u>s</u>		28	yuppy	103
}		28	yuppy	104
<b>reset</b> s to mark		31	lubber	101
advance r				
}				
}				

<pre>while not done {</pre>
while (r < s) { advance r }
while $(r > s)$ { advance s }
<pre>mark s // save start of "block"</pre>
while (r == s) {
// Outer loop over r
while $(r == s)$ {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
<b>reset</b> s to mark
advance r
}

sid	sname		sid	bid
22	dustin		28	103
28	уирру		28	104
31	lubbe	r	31	101
31	lubber2		<b>&gt;</b> 31	102
44	guppy		42	142
57	rusty		58	107
s 2	sid 28	sname yuppy	bid 103	
2	28	yuppy	104	
3	31	lubber	101	

<pre>while not done {</pre>	sid
while $(r < s) \{ advance r \}$	22
while (r > s) { advance s }	28
<pre>mark s // save start of "block"</pre>	━→31
while $(r == s)$ {	31
// Outer loop over r	44
while $(r == s)$ {	57
// Inner loop over s	
yield <r, s=""></r,>	
advance s	
}	
<b>reset</b> s to mark	
advance r	
}	

}

sid	snar	ne	sid	bid
22	dust	in	28	103
28	yupp	ру	28	104
31	lubb	er	31	101
31	lubb	er2	<b>=</b> 31	102
44	gup	су	42	142
57	rust	ý	58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	

<pre>while not done {</pre>
while $(r < s)$ { advance r }
while $(r > s)$ { advance s }
mapping $(/ contact of (b) contact)$
mark S // save start of "DIOCK"
while (r == s) {
// Outer loop over r
while $(r == s)$ {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
<b>reset</b> s to mark
duvance r
}

- 1 - 1			aid	ام زما
SIC	snam	е	SIQ	DIQ
22	dustir	า	28	103
28	yuppy	/	28	104
31	lubbe	r	31	101
31	lubbe	r2	31	102
44	guppy	/	<b>4</b> 2	142
57	rusty		58	107
	sid	sname	bid	
28		yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	

<pre>while not done {</pre>
while (r < s) { advance r }
while $(r > s)$ { advance s }
<pre>mark s // save start of "block"</pre>
while $(r = s)$ {
// Outen loop oven n
// Outer toop over th
while (r == s) {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
<b>reset</b> s to mark
advance r
}
}

sid	snar	ne	sid	bid
22	dust	in	28	103
28	yupp	ру	28	104
31	lubb	er	<b>&gt;</b> 31	101
31	lubb	er2	31	102
44	gupp	ру	42	142
57	rusty	/	58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	

<pre>while not done {</pre>
while $(r < s)$ { advance r }
while $(r > s)$ { advance s }
<pre>mark s // save start of "block"</pre>
while (r == s) {
// Outer loop over r
while (r == s) {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
reset s to mark
advance r
}

sid	sna	me	sid	bid
22	dust	tin	28	103
28	yup	ру	28	104
31	lubb	ber	<b>→</b> 31	101
⇒31	lubb	per2	31	102
44	gup	ру	42	142
57	rust	У	58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	

<pre>while not done {    while (r &lt; s) { advance r }</pre>	
while (r > s) { advance s }	
<pre>mark s // save start of "block"</pre>	
while (r == s) {	
// Outer loop over r	
while (r == s) {	
// Inner loop over s	
yield <r, s=""></r,>	
advance s	
}	
<b>reset</b> s to mark	
advance r	
}	
}	

sid	sna	me	sid	bid
22	dust	tin	28	103
28	yup	ру	28	104
31	lubb	ber	<b>&gt;</b> 31	101
♦31	lubb	per2	31	102
44	gup	ру	42	142
57	rust	у	58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	

<pre>while not done {</pre>
while (r < s) { advance r }
while (r > s) {
<pre>mark s // save start of "block"</pre>
while (r == <u>s</u> ) {
// Outer loop over r
while (r == <mark>s</mark> ) {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
<b>reset</b> s to mark
advance r
}
} ´

sid	sna	me	sid	bid
22	dus	tin	28	103
28	yup	ру	28	104
31	lubl	ber	31	101
⇒31	lubl	per2	<b>=</b> 31	102
44	gup	ру	42	142
57	rust	.y	58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	

<pre>while not done {</pre>
while $(r < s)$ { advance r }
while $(r > s)$ { advance s }
Mark S // save start of "block"
while (r == s) {
// Outer loop over r
while (r == s) {
// Inner loop over s
yield <r, s=""></r,>
advance
}
<b>reset</b> s to mark
advance r
}
}

sid	snam	е	sid	bid
22	dusti	า	28	103
28	yuppy	Y	28	104
31	lubbe	er	31	101
31	lubbe	er2	<b>=</b> 31	102
44	guppy	ý	42	142
57	rusty		58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	
	31	lubber2	102	

sid	snam	e	sid	bid
22	dustir	ı	28	103
28	yuppy	y	28	104
31	lubbe	r	31	101
31	lubbe	r2	31	102
44	guppy	/	➡ 42	142
57	rusty		58	107
			6.1.1	
	SIC	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	
	31	lubber2	102	

<pre>while not done {</pre>
while $(r < s)$ { advance r }
while $(n > c) \\ \begin{cases} advance \\ c \\ \end{cases}$
mark s // save start of "block"
while (r == <u>s</u> ) {
// Outer loop over r
while (r == s) {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
reset s to mark
advance
}
}
•

sid	snam	e	sid	bid
22	dusti	n	28	103
28	yupp	y	28	104
31	lubbe	er	<b>&gt;</b> 31	101
31	lubbe	er2	31	102
44	gupp	у	42	142
57	rusty		58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	
	31	lubber2	102	

while not done {
while $(r < s)$ { advance r }
while $(r > s)$ { advance s }
mark s // save start of "block"
while (r == s) {
// Outer loop over r
while (r == <u>s</u> ) {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
reset s to mark
advance
}
}

sid	sna	me	sid	bid
22	dus	tin	28	103
28	yup	ру	28	104
31	lub	ber	<b></b> 31	101
31	lub	per2	31	102
•44	gup	ру	42	142
57	rust	ý	58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	
	31	lubber2	102	

<pre>while not done {</pre>
while $(r < s)$ { advance r }
while $(r > s) \{ advance s \}$
<pre>mark s // save start of "block"</pre>
while (r == s) {
// Outer loop over r
while (r == s) {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
<b>reset</b> s to mark
advance r
}
}

sid	snam	е	sid	bid
22	dustir	ı	28	103
28	yuppy	/	28	104
31	lubbe	r	<b>&gt;</b> 31	101
31	lubbe	r2	31	102
44	guppy	/	42	142
57	rusty		58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	
	31	lubber2	102	

<pre>while not done {</pre>		
<pre>while (r &lt; s) { advance r }</pre>		
while $(r > s)$ { advance s }		
<pre>mark s // save start of "block"</pre>		
while (r == s) {		
// Outer loop over r		
while $(r == s)$ {		
// Inner loop over s		
yield <r, s=""></r,>		
advance s		
}		
<b>reset</b> s to mark		
advance r		
}		
}		

sid	snam	е	sid	bid
22	dusti	n	28	103
28	yuppy	y	28	104
31	lubbe	er	31	101
31	lubbe	er2	<b>=</b> 31	102
44	gupp	y	42	142
57	rusty		58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	
	31	lubber2	102	

<pre>while not done {</pre>	S
<pre>while (r &lt; s) { advance r }</pre>	2
<pre>while (r &gt; s) { advance s }</pre>	2
mark s // save start of "block"	3
while $(r == s)$ {	3
// Outer loop over r	-4
while (r == s) {	5
// Inner loop over s	
yield <r, s=""></r,>	
advance s	
}	
<b>reset</b> s to mark	
advance r	
}	
}	

sid	snai	me	sid	bid
22	dust	tin	28	103
28	yup	ру	28	104
31	lubb	ber	31	101
31	lubb	per2	31	102
44	gup	ру	<b>4</b> 2	142
57	rust	у	58	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	
	31	lubber2	102	

<pre>while not done {</pre>
<pre>while (r &lt; s) { advance r }</pre>
<pre>while (r &gt; s) { advance s }</pre>
<pre>mark s // save start of "block"</pre>
while (r == s) {
// Outer loop over r
while $(r == s)$ {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
reset s to mark
advance r
}
}

sid	snam	e	sid	bid
22	dustir	ı	28	103
28	yuppy	/	28	104
31	lubbe	r	31	101
31	lubbe	r2	31	102
44	guppy	/	42	142
57	rusty		<b>→</b> 58	107
:	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
:	31	lubber	102	
	31	lubber2	101	

<pre>while not done {</pre>
<pre>while (r &lt; s) { advance r }</pre>
<pre>while (r &gt; s) { advance s }</pre>
<pre>mark s // save start of "block"</pre>
while (r == s) {
// Outer loop over r
while $(r == s)$ {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
reset s to mark
advance r
}
}

sid	snam	e	sid	bid
22	dustir	ı	28	103
28	yuppy	/	28	104
31	lubbe	r	31	101
31	lubbe	r2	31	102
44	guppy	/	42	142
57	rusty		<b>→</b> 58	107
:	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
	31	lubber	101	
:	31	lubber	102	
	31	lubber2	101	

<pre>while not done {    while (r &lt; s) { advance r }    while (r &gt; s) { advance s }</pre>	
<pre>mark s // save start of "block"</pre>	
while (r == s) {	
// Outer loop over r	
while (r == <u>s</u> ) {	
// Inner loop over s	
yield <r, s=""></r,>	
advance s	
}	
reset s to mark	
advance r	
} }	

sid	snam	е	sid	bid
22	dustir	ı	28	103
28	yuppy	/	28	104
31	lubbe	r	31	101
31	lubbe	r2	31	102
44	guppy	/	42	142
57	rusty		<del></del>	107
	sid	sname	bid	
	28	yuppy	103	
	28	yuppy	104	
:	31	lubber	101	
	31	lubber	102	
	31	lubber2	101	

<pre>while not done {</pre>
while $(r < s)$ { advance r }
while $(r > s)$ { advance s }
<pre>mark s // save start of "block"</pre>
while ( <u>r == s</u> ) {
// Outer loop over r
while (r == s) {
// Inner loop over s
yield <r, s=""></r,>
advance s
}
<b>reset</b> s to mark
advance r
}
}

	sid	snam	е	sid	bid
	22	dustir	ı	28	103
	28	уирру		28	104
	31	lubber		31	101
	31	lubber2		31	102
	44	guppy		42	142
	57	rusty		<del></del>	107
	c	sid	sname	bid	
	· · · ·	510	Shame	bid	
		28	yuppy	103	
		28	yuppy	104	
		31	lubber	101	
		31	lubber	102	
		31	lubber2	101	

bid

<pre>while not done {</pre>	sid	snai	ne	sid
<pre>while (r &lt; s) { advance r }</pre>	22	dust	in	28
<pre>while (r &gt; s) { advance s }</pre>	28	yup	ру	28
mark s // save start of "block"	31	lubb	ber	31
while $(r == s)$ {	31	lubb	per2	31
// Outer loop over r	44	gup	ру	42
while $(r == s)$ {	57	rust	у	➡58
// Inner loop over s 🛛 💻	•			
yield <r, s=""></r,>		sid	sname	bid
advance s		28	yuppy	103
}		28	yuppy	104
<b>reset</b> s to mark		31	lubber	101
advance r		31	lubber	102
} \		31	lubber2	101
J		31	lubber2	102

### Sort-Merge Join (SMJ)

- I/O cost?
  - Cost of sorting R
  - Cost of sorting S
  - The merge step: [R] + [S]
    - Only one pass (if we assume there aren't a lot of duplicates)

### Sort-Merge Join (SMJ)

- An optimization we can sometimes make
  - Recall materialization (write + read) is expensive
  - We only have to (assuming no duplicate values in R) make one pass through the sorted relation → we don't need the sorted relations to be materialized!
  - In the final merge pass of sorting both relations, instead of writing the sorted relations to disk, we can stream them into the second part of SMJ!
    - Reduces I/O cost by 2\*([R] + [S])!

### Sort-Merge Join (SMJ)

- An optimization we can sometimes make
  - In the final merge pass of sorting both relations, instead of writing the sorted relations to disk, we can stream them into the second part of SMJ!
    - Since we are iterating over R and S anyway, we can begin outputting what will join from the two relations
    - We have to be able to fit the input buffers of the last merge pass of sorting R and sorting S in memory, as well as have one output buffer for joined tuples
    - Need: # runs in last merge pass for R + # runs in last merge pass for S ≤ B - 1

### Worksheet Q1e

In the average case, how many disk I/Os are needed to perform a sort-merge join (unoptimized/optimized)?

Companies: (company\_id, industry, ipo\_date) Nyse: (company\_id, date, trade, quantity)

- 20 pages of memory
- [N] = 100 pages, p<sub>N</sub> = 100 tuples per page
- [C] = 50 pages,  $p_c = 50$  tuples per page

### Worksheet Q1e

How many disk I/Os are needed to perform a sort-merge join (unoptimized/optimized)?

### **Unoptimized:**

### Sorting N:

```
Pass 1 - ceil(100/20) = 5 sorted runs of 20 pages
```

each

Pass 2 - ceil(5/19) = 1 sorted run of 100 pages each

Total I/Os: 4 \* (100 pages) = 400 I/Os

Companies: (company\_id, industry, ipo\_date)
Nyse: (company\_id, date, trade, quantity)

- 20 pages of memory
- [N] = 100 pages, p<sub>N</sub> = 100 tuples per page

• [C] = 50 pages, 
$$p_C = 50$$
 tuples per page

### Sorting C:

**Pass 1** - ceil(50/20) = 3 sorted runs of 20

pages, 20 pages, and 10 pages

Pass 2 - ceil(3/19) = 1 sorted run of 50 pages

Total I/Os: 4 \* (50 pages) = 200 I/Os

**Merging:** [C] + [N] = 150 I/Os

Total SMJ I/Os: 200 + 400 + 150 = 750 I/Os

#### Unoptimized SMJ



#### I/Os:

- Sorting N: 400
- Sorting C: ??
- Merging: ??

Read, Write, Merge

#### Unoptimized SMJ



#### I/Os:



#### • Sorting C: 200

• Merging: ??

Read, Write, Merge

#### Sorted Relation N

1 sorted run of 100 pages

Unoptimized SMJ



I/Os: 750

- Sorting N: 400
- Sorting C: 200
- Merging: 150

Read, Write, Merge

### Worksheet Q1e

How many disk I/Os are needed to perform a sort-merge join (unoptimized/optimized)?

Can we perform the SMJ optimization?

Sorting N:

**Pass 0** - ceil(100/20) = 5 sorted runs of 20 pages each

Pass 1 - ceil(5/19) = 1 sorted run of 100 pages each

Total I/Os: 4 \* (100 pages) = 400 I/Os

Companies: (company\_id, industry, ipo\_date)
Nyse: (company\_id, date, trade, quantity)

- 20 pages of memory
- [N] = 100 pages, p<sub>N</sub> = 100 tuples per page

• [C] = 50 pages, 
$$p_c = 50$$
 tuples per page

Sorting C: Pass 0 - ceil(50/20) = 3 sorted runs of 20 pages, 20 pages, and 10 pages Pass 1 - ceil(3/19) = 1 sorted run of 50 pages Total I/Os: 4 \* (50 pages) = 200 I/Os

### Worksheet Q1e

How many disk I/Os are needed to perform a sort-merge join (unoptimized/optimized)?

### Can we perform the SMJ optimization?

Yes.

During the 2nd to last pass, we produce 5 sorted runs of N and 3 sorted runs of C. Since the number of runs of C + the number of runs of N  $\leq$  20 - 1, we can optimize sort merge join and combine the last sorting pass and final merging pass to save 2 \* ([C] + [N]) I/Os. Total I/Os = 750 - 2(50+100) = 450 I/Os

Companies: (company\_id, industry, ipo\_date) Nyse: (company\_id, date, trade, quantity)

- 20 pages of memory
- [N] = 100 pages,  $p_N = 100$  tuples per page
- [C] = 50 pages,  $p_C = 50$  tuples per page

### **Optimized SMJ**



• i.e. saves 2([C] + [N]) I/Os

Read, Write

I/Os: cost(unoptimized SMJ) - 2([C] + [N]) = 750 - 2(150) = 450

### **Grace Hash Join**

- Same idea as SMJ, but let's build some hash tables instead
- Two passes: **partition** the data, then **build** an in-memory hash table and **probe** it
  - First, partition R and S into B 1 partitions (like in external hashing), *using the same hash function* 
    - All the tuples in R matching a tuple in S must be in the same partition → we can consider each partition independently

### Grace Hash Join

- Same idea as SMJ, but let's build some hash tables instead
- Two passes: **partition** the data, then **build** an in-memory hash table and **probe** it
  - Then, build an in-memory hash table for a partition of R
  - We can use this in-memory hash table to find all the tuples in R that match a tuple in S
    - Stream in tuples of S, probe the hash table, output matching tuples

### Grace Hash Join: Partition



### Grace Hash Join: Partition



### Grace Hash Join: Partition














R S







S

R



S

R



#### **Grace Hash Join**

- We need partitions of R (but not S) to fit in B 2 pages
  - 1 page reserved for streaming S partition
  - 1 page reserved for streaming output
- What if partitions of R are too big?
  - If S is smaller, do S  $\bowtie_{\theta}$  R instead
  - Recursively partition! Make sure that for any partition of R you recursively partition, the matching S partition is also recursively partitioned!

#### **Grace Hash Join**

Pass 2: Build and Probe

- Build an in-memory hash table for a partition of R
- Stream in tuples of S, probe the hash table, output matching tuples



Build in-memory hash table of R and stream in tuples of S



Build in-memory hash table of R and stream in tuples of S











Build in-memory hash table of R and stream in tuples of S



Build in-memory hash table of R and stream in tuples of S











. . .

#### Worksheet

#### Worksheet Q2a

If we had 10 buffer pages, how many partitioning phases would we require for grace hash join?

- 2 tables: Catalog and Transactions
- [C] = 100 pages,  $p_c = 20$  tuples per page
- [T] = 50 pages,  $p_T = 50$  tuples per page
- Assume the hash functions uniformly distribute the data for both tables.

#### Worksheet Q2a

If we had 10 buffer pages, how many partitioning phases would we require for grace hash join?

T is smaller, so we need its partitions to be at most B - 2 = 8 pages. After 1 partitioning pass, we have partitions of size 6, which is <= 8 so we only need **1 partitioning pass**.

- 2 tables: Catalog and Transactions
- [C] = 100 pages,  $p_c = 20$  tuples per page
- [T] = 50 pages,  $p_T = 50$  tuples per page
- Assume the hash functions uniformly distribute the data for both tables.

#### Worksheet Q2b

What is the IO cost for the grace hash join then? Assume uniform partitioning.

- 2 tables: Catalog and Transactions
- [C] = 100 pages,  $p_{C}$  = 20 tuples per page
- [T] = 50 pages,  $p_T = 50$  tuples per page
- Assume the hash functions uniformly distribute the data for both tables.

#### Worksheet Q2b

What is the IO cost for the grace hash join then? Assume uniform partitioning.

- 2 tables: Catalog and Transactions
- [C] = 100 pages,  $p_C = 20$  tuples per page
- [T] = 50 pages,  $p_T = 50$  tuples per page
- Assume the hash functions uniformly distribute the data for both tables.

We need 1 partitioning pass.

#### **Partitioning phase:**

ceil([C]/(B - 1)) = 12 pages per partition for C, 12(9) pages in total after partitioning

ceil([T]/(B - 1)) = 6 pages per partition for T, 6(9) pages in total after partitioning

**Partitioning IOs:** 100 I/Os to read from Catalog + 12(9) to write for Catalog + 50 I/Os to read from Transactions + 6(9) to write for Transactions = 312 I/Os

Probing phase: 12(9) + 6(9) = 162 I/Os to read from Catalog and Transactions

Total: 312 + 162 = 474 I/Os

#### Worksheet Q2c

If we only had 8 buffer pages, how many partitioning phases would there be?

- 2 tables: Catalog and Transactions
- [C] = 100 pages,  $p_c = 20$  tuples per page
- [T] = 50 pages,  $p_T = 50$  tuples per page
- Assume the hash functions uniformly distribute the data for both tables.

#### Worksheet Q2c,d

- B = 8
  - [<u>C]</u> = 100 pages
  - [T] = 50 pages
- Assume the hash functions uniformly distribute the data for both tables.



fit in B-2 = 6 pages, so we

must recursively partition.

fit in B-2 = 6 pages, so we can enter the Build and Probe phase.

#### Worksheet Q2c,d

- B = 8
- = 100 pages = 50 pages
- Assume the hash functions uniformly distribute the data for both tables.



#### Worksheet Q2c,d

- B = 8
- = 100 pages
  - = 50 pages
- Assume the hash functions uniformly distribute the data for both tables.



Note: We can alternatively build a hash table on T and probe C since partitions for either relation fit in B-2 pages.

#### Worksheet Q2c

If we only had 8 buffer pages, how many partitioning phases would there be?

T is smaller, so we need its partitions to be at most B -

2 = 6 pages. After 1 partitioning pass,

we have partitions of size 8, which is too big to fit in

B-2 buffer pages. We need a second

partitioning pass.  $8 / 7 = 1.1 \rightarrow 2$  pages, which is small

enough to fit in B-2 buffer pages.

Therefore, we need **2 passes** in total.

- 2 tables: Catalog and Transactions
- [C] = 100 pages,  $p_c = 20$  tuples per page
- [T] = 50 pages,  $p_T = 50$  tuples per page
- Assume the hash functions uniformly distribute the data for both tables.

#### Worksheet Q2d

What will be the IO cost?

- 2 tables: Catalog and Transactions
- [C] = 100 pages,  $p_C = 20$  tuples per page
- [T] = 50 pages,  $p_T = 50$  tuples per page
- Assume the hash functions uniformly distribute the data for both tables.
## Worksheet Q2d

What will be the IO cost?

Partitioning phase:

ceil([C]/(B - 1)) = 15 pages per partition for C

ceil([T]/(B - 1)) = 8 pages per partition for T

ceil([C]/(B - 1)) = 3 pages per partition for second pass for C

ceil([T]/(B - 1)) = 2 pages per partition for second pass for T

Read 1st Write 1st Read 2nd Write 2nd

Partitioning IOs: [100 + 50] + [15(7) + 8(7)] + [15(7) + 8(7)] + [3(49) + 2(49)] = 717 I/Os

Build and Probe Phase: 3(49) + 2(49) = 245 IOs

Total: 717 + 245 = **962 I/Os** 

- 2 tables: Catalog and Transactions
- [C] = 100 pages,  $p_c = 20$  tuples per page
- [T] = 50 pages,  $p_T = 50$  tuples per page
- Assume the hash functions uniformly distribute the data for both tables.

## Attendance Link

https://cs186berkeley.net/attendance

