

Database System Internals Operator Algorithms (part 2)

Paul G. Allen School of Computer Science and Engineering
University of Washington, Seattle

January 27, 2020 CSE 444 - Winter 2020 1

1

Announcements

- Homework 2 released
 - Due January 31st
- 544 paper 1 report due Today
- Lab 2 posted after class
 - Part 1 (operator algos) due Friday
 - Part 2 (insert/delete support) due following Friday

January 27, 2020 CSE 444 - Winter 2020 2

2

Block-Memory Refinement

```

for each group of M-1 pages r in R do
  for each page of tuples s in S do
    for all pairs of tuples t1 in r, t2 in s
      if t1 and t2 join then output (t1, t2)
        
```

What is the Cost?

January 27, 2020 CSE 444 - Winter 2020 3

3

Block Memory Refinement

M=3

Disk

Patient		Insurance	
1	2	2	4
3	4	4	3
9	6	2	8
8	5	8	9

Input buffer for Patient

Input buffer for Insurance

No output buffer: stream to output

4

Block Memory Refinement

M=3

Disk

Patient		Insurance	
1	2	2	4
3	4	4	3
9	6	2	8
8	5	8	9

Input buffer for Patient

Input buffer for Insurance

No output buffer: stream to output

5

Block Memory Refinement

M=3

Disk

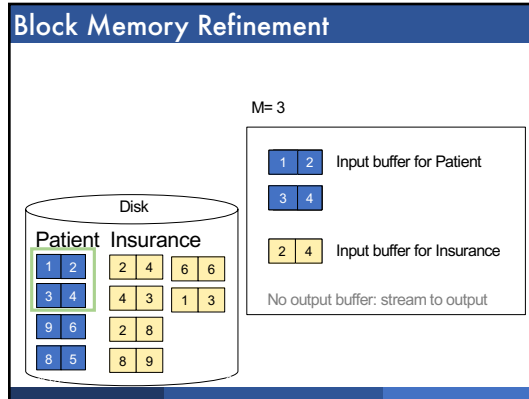
Patient		Insurance	
1	2	2	4
3	4	4	3
9	6	2	8
8	5	8	9

Input buffer for Patient

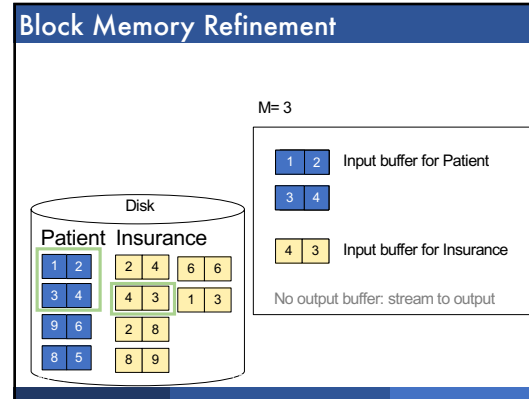
Input buffer for Insurance

No output buffer: stream to output

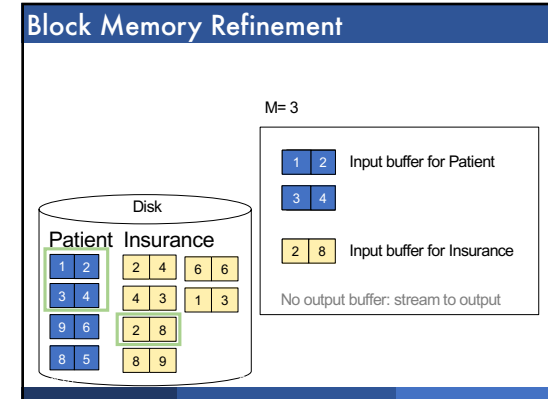
6



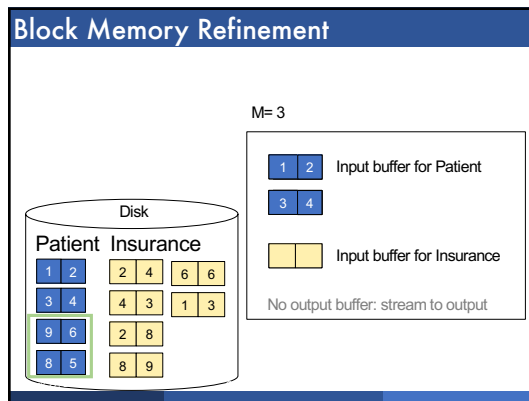
7



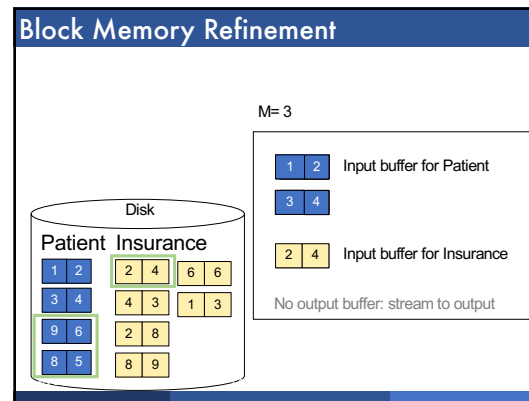
8



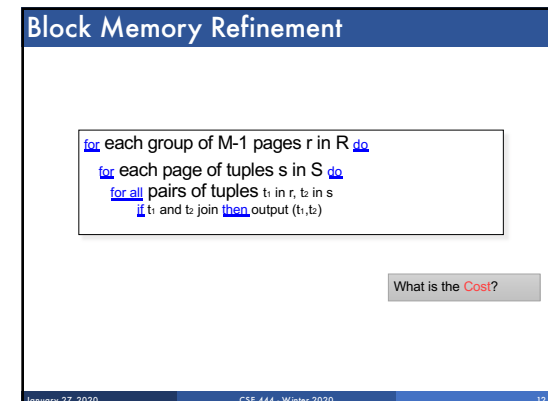
9



10



11



12

Block Memory Refinement

```

for each group of M-1 pages r in R do
  for each page of tuples s in S do
    for all pairs of tuples t1 in r, t2 in s
      if t1 and t2 join then output (t1, t2)
  
```

Cost: $B(R) + B(R)B(S)/(M-1)$

What is the Cost?

13

Outline

- Join operator algorithms
 - One-pass algorithms (Sec. 15.2 and 15.3)
 - Index-based algorithms (Sec 15.6)
 - Two-pass algorithms (Sec 15.4 and 15.5)

19

Index Based Selection

Selection on equality: $\sigma_{a=v}(R)$

- $B(R)$ = size of R in blocks
- $T(R)$ = number of tuples in R
- $V(R, a)$ = # of distinct values of attribute a

20

Index Based Selection

Selection on equality: $\sigma_{a=v}(R)$

- $B(R)$ = size of R in blocks
- $T(R)$ = number of tuples in R
- $V(R, a)$ = # of distinct values of attribute a

What is the cost in each case?

- Clustered index on a:
- Unclustered index on a:

21

Index Based Selection

Selection on equality: $\sigma_{a=v}(R)$

- $B(R)$ = size of R in blocks
- $T(R)$ = number of tuples in R
- $V(R, a)$ = # of distinct values of attribute a

What is the cost in each case?

- Clustered index on a: $B(R)/V(R, a)$
- Unclustered index on a:

22

Index Based Selection

Selection on equality: $\sigma_{a=v}(R)$

- $B(R)$ = size of R in blocks
- $T(R)$ = number of tuples in R
- $V(R, a)$ = # of distinct values of attribute a

What is the cost in each case?

- Clustered index on a: $B(R)/V(R, a)$
- Unclustered index on a: $T(R)/V(R, a)$

23

Index Based Selection

Selection on equality: $\sigma_{a=v}(R)$

- $B(R)$ = size of R in blocks
- $T(R)$ = number of tuples in R
- $V(R, a)$ = # of distinct values of attribute a

What is the cost in each case?

- Clustered index on a : $B(R)/V(R, a)$
- Unclustered index on a : $T(R)/V(R, a)$

Note: we ignore I/O cost for index pages

January 27, 2020 CSE 444 - Winter 2020 24

24

Index Based Selection

▪ Example:

$B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$

cost of $\sigma_{a=v}(R) = ?$

- Table scan:
- Index based selection:

January 27, 2020 CSE 444 - Winter 2020 25

25

Index Based Selection

▪ Example:

$B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$

cost of $\sigma_{a=v}(R) = ?$

- Table scan: $B(R) = 2,000$ I/Os
- Index based selection:

January 27, 2020 CSE 444 - Winter 2020 26

26

Index Based Selection

▪ Example:

$B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$

cost of $\sigma_{a=v}(R) = ?$

- Table scan: $B(R) = 2,000$ I/Os
- Index based selection:
 - If index is clustered:
 - If index is unclustered:

January 27, 2020 CSE 444 - Winter 2020 27

27

Index Based Selection

▪ Example:

$B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$

cost of $\sigma_{a=v}(R) = ?$

- Table scan: $B(R) = 2,000$ I/Os
- Index based selection:
 - If index is clustered: $B(R)/V(R, a) = 100$ I/Os
 - If index is unclustered:

January 27, 2020 CSE 444 - Winter 2020 28

28

Index Based Selection

▪ Example:

$B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$

cost of $\sigma_{a=v}(R) = ?$

- Table scan: $B(R) = 2,000$ I/Os
- Index based selection:
 - If index is clustered: $B(R)/V(R, a) = 100$ I/Os
 - If index is unclustered: $T(R)/V(R, a) = 5,000$ I/Os

January 27, 2020 CSE 444 - Winter 2020 29

29

Index Based Selection

Example:

$B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$

cost of $\sigma_{a=v}(R) = ?$

- Table scan: $B(R) = 2,000$ I/Os!
- Index based selection:
 - If index is clustered: $B(R)/V(R,a) = 100$ I/Os
 - If index is unclustered: $T(R)/V(R,a) = 5,000$ I/Os!

January 27, 2020 CSE 444 - Winter 2020 30

30

Index Based Selection

Example:

$B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$

cost of $\sigma_{a=v}(R) = ?$

- Table scan: $B(R) = 2,000$ I/Os!
- Index based selection:
 - If index is clustered: $B(R)/V(R,a) = 100$ I/Os
 - If index is unclustered: $T(R)/V(R,a) = 5,000$ I/Os!

Lesson: Don't build unclustered indexes when $V(R,a)$ is small !

January 27, 2020 CSE 444 - Winter 2020 31

31

Index Based Selection

Example:

$B(R) = 2000$
 $T(R) = 100,000$
 $V(R, a) = 20$

cost of $\sigma_{a=v}(R) = ?$

- Table scan: $B(R) = 2,000$ I/Os
- Index based selection:
 - If index is clustered: $B(R)/V(R,a) = 100$ I/Os
 - If index is unclustered: $T(R)/V(R,a) = 5,000$ I/Os

Lesson: Don't build unclustered indexes when $V(R,a)$ is small !

January 27, 2020 CSE 444 - Winter 2020 32

32

Index Nested Loop Join

$R \bowtie S$

- Assume S has an index on the join attribute
- Iterate over R , for each tuple fetch corresponding tuple(s) from S
- Previous nested loop join: cost
 - $B(R) + T(R) * B(S)$
- Index Nested Loop Join Cost:**
 - If index on S is clustered: $B(R) + T(R)B(S)/V(S,a)$
 - If index on S is unclustered: $B(R) + T(R)T(S)/V(S,a)$

January 27, 2020 CSE 444 - Winter 2020 33

33

Outline

- Join operator algorithms**
 - One-pass algorithms (Sec. 15.2 and 15.3)
 - Index-based algorithms (Sec 15.6)
 - Two-pass algorithms (Sec 15.4 and 15.5)

January 27, 2020 CSE 444 - Winter 2020 34

34

Two-Pass Algorithms

- Fastest algorithm seen so far is one-pass hash join
 What if data does not fit in memory?
- Need to process it in multiple passes
- Two key techniques
 - Sorting
 - Hashing

January 27, 2020 CSE 444 - Winter 2020 35

35

Basic Terminology

- A run in a sequence is an increasing subsequence
- What are the runs?
2, 4, 99, 103, 88, 77, 3, 79, 100, 2, 50

January 27, 2020 CSE 444 - Winter 2020 36

36

Basic Terminology

- A run in a sequence is an increasing subsequence
- What are the runs?
2, 4, 99, 103, 88, 77, 3, 79, 100, 2, 50

January 27, 2020 CSE 444 - Winter 2020 37

37

External Merge-Sort: Step 1

Phase one: load M blocks in memory, sort, send to disk, repeat

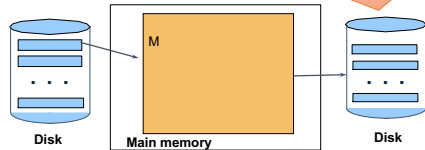
January 27, 2020 CSE 444 - Winter 2020 41

41

External Merge-Sort: Step 1

Phase one: load M blocks in memory, sort, send to disk, repeat

Q: How long are the runs?



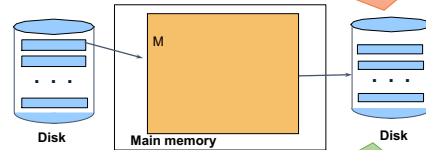
January 27, 2020 CSE 444 - Winter 2020 42

42

External Merge-Sort: Step 1

Phase one: load M blocks in memory, sort, send to disk, repeat

Q: How long are the runs?



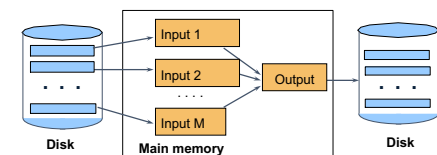
January 27, 2020 CSE 444 - Winter 2020 43

43

A: Length = M blocks

Phase two: merge M runs into a bigger run

- Merge M – 1 runs into a new run
- Result: runs of length M (M – 1) $\approx M^2$



January 27, 2020 CSE 444 - Winter 2020 45

45

Example

- Merging three runs to produce a longer run:

0, 14, 33, 88, 92, 192, 322
 2, 4, 7, 43, 78, 103, 523
 1, 6, 9, 12, 33, 52, 88, 320

Output:
 0

46

Example

- Merging three runs to produce a longer run:

0, 14, 33, 88, 92, 192, 322
 2, 4, 7, 43, 78, 103, 523
 1, 6, 9, 12, 33, 52, 88, 320

Output:
 0, ?

47

Example

- Merging three runs to produce a longer run:

0, 14, 33, 88, 92, 192, 322
 2, 4, 7, 43, 78, 103, 523
 1, 6, 9, 12, 33, 52, 88, 320

Output:
 0, 1, ?

48

Example

- Merging three runs to produce a longer run:

0, 14, 33, 88, 92, 192, 322
 2, 4, 7, 43, 78, 103, 523
 1, 6, 9, 12, 33, 52, 88, 320

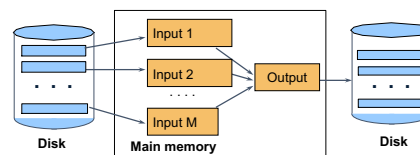
Output:
 0, 1, 2, 4, 6, 7, ?

49

External Merge-Sort: Step 2

Phase two: merge M runs into a bigger run

- Merge M - 1 runs into a new run
- Result: runs of length M (M - 1) \approx M²



If approx. $B \leq M^2$ then we are done

50

Cost of External Merge Sort

- Assumption: $B(R) \leq M^2$
- Read+write+read = $3B(R)$

51

Discussion

- What does $B(R) \leq M^2$ mean?
- How large can R be?

January 27, 2020 CSE 444 - Winter 2020 52

52

Discussion

- What does $B(R) \leq M^2$ mean?
- How large can R be?
- Example:
 - Page size = 32KB
 - Memory size 32GB: $M = 10^6$ pages

January 27, 2020 CSE 444 - Winter 2020 53

53

Discussion

- What does $B(R) \leq M^2$ mean?
- How large can R be?
- Example:
 - Page size = 32KB
 - Memory size 32GB: $M = 10^6$ pages
- R can be as large as 10^{12} pages
 - 32×10^{15} Bytes = 32 PB

January 27, 2020 CSE 444 - Winter 2020 54

54

Merge-Join

- Join $R \bowtie S$
- How?....

January 27, 2020 CSE 444 - Winter 2020 55

55

Merge-Join

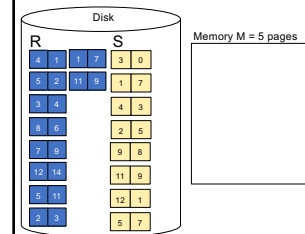
- Join $R \bowtie S$
- Step 1a: generate initial runs for R
 - Step 1b: generate initial runs for S
 - Step 2: merge and join
 - Either merge first and then join
 - Or merge & join at the same time

January 27, 2020 CSE 444 - Winter 2020 56

56

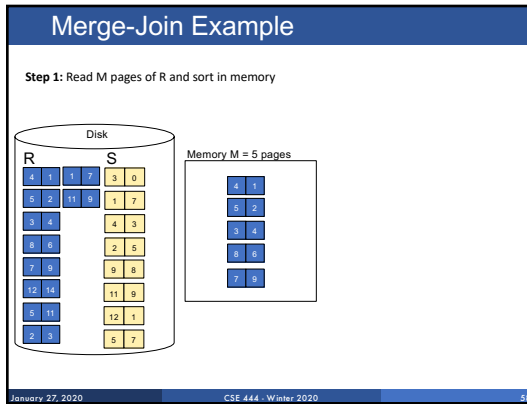
Merge-Join Example

Setup: Want to join R and S
 Relation R has 10 pages with 2 tuples per page
 Relation S has 8 pages with 2 tuples per page
 Values shown are values of join attribute for each given tuple

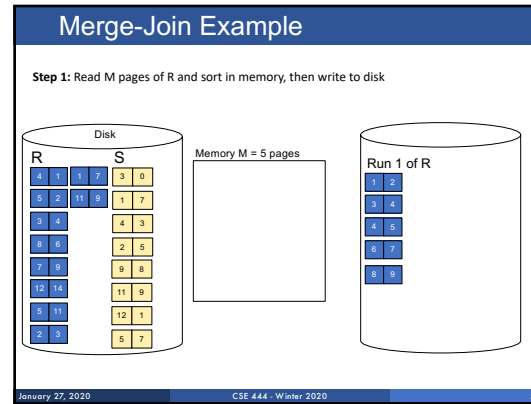


January 27, 2020 CSE 444 - Winter 2020 57

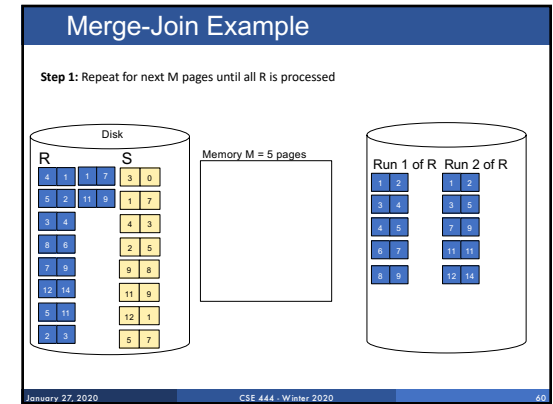
57



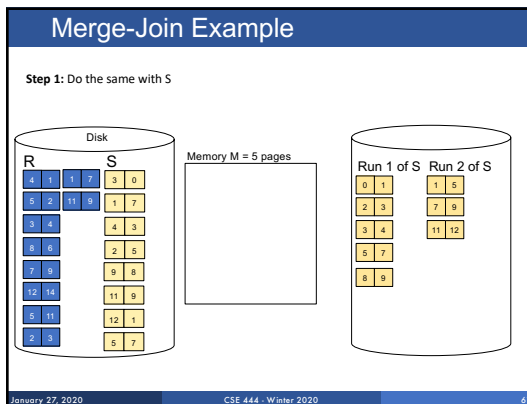
58



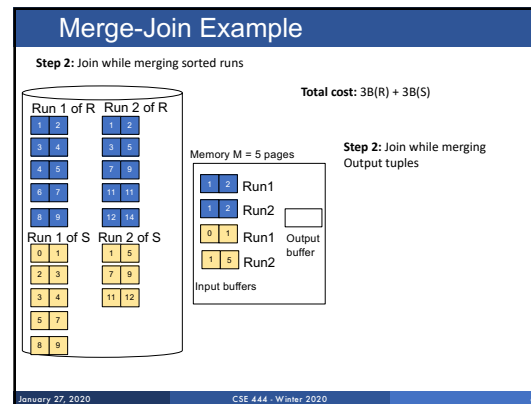
59



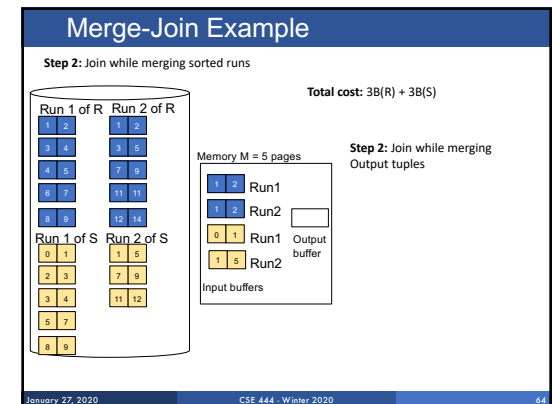
60



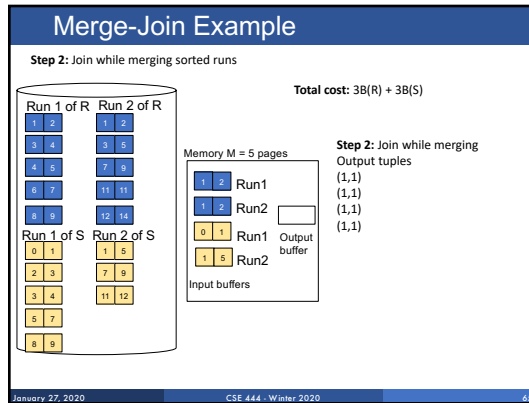
61



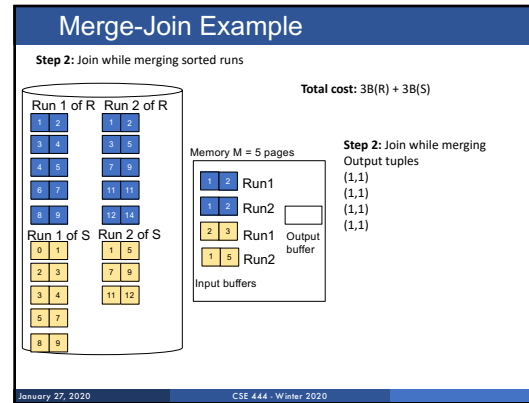
62



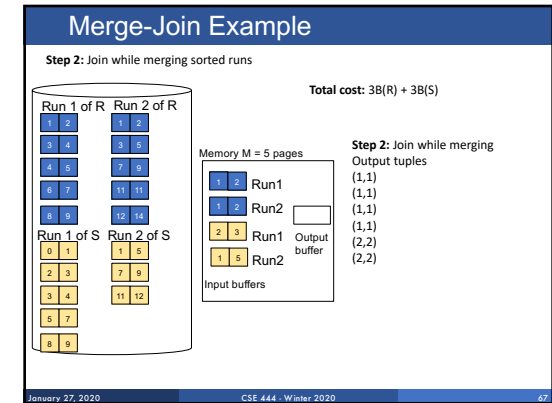
64



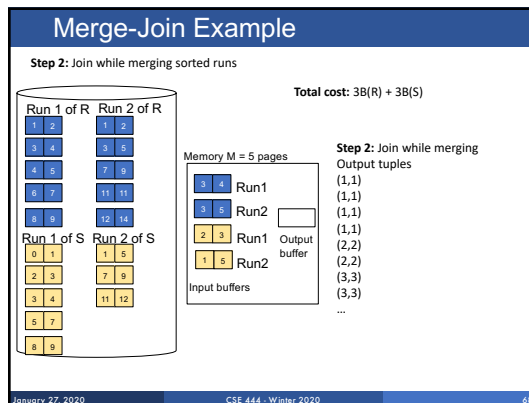
65



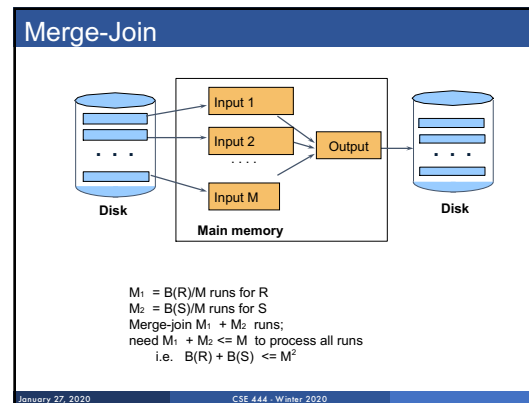
66



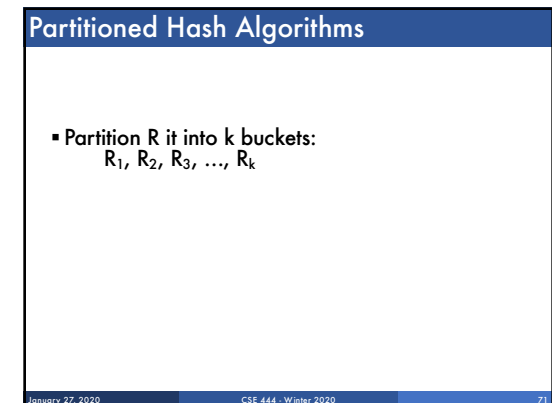
67



68



70



71

Partitioned Hash Algorithms

- Partition R it into k buckets:
 $R_1, R_2, R_3, \dots, R_k$
- Assuming $B(R_1)=B(R_2)=\dots=B(R_k)$, we have
 $B(R_i) = B(R)/k$, for all i

January 27, 2020 CSE 444 - Winter 2020 72

72

Partitioned Hash Algorithms

- Partition R it into k buckets:
 $R_1, R_2, R_3, \dots, R_k$
- Assuming $B(R_1)=B(R_2)=\dots=B(R_k)$, we have
 $B(R_i) = B(R)/k$, for all i
- Goal: each R_i should fit in main memory:
 $B(R_i) \leq M$

January 27, 2020 CSE 444 - Winter 2020 73

73

Partitioned Hash Algorithms

- Partition R it into k buckets:
 $R_1, R_2, R_3, \dots, R_k$
- Assuming $B(R_1)=B(R_2)=\dots=B(R_k)$, we have
 $B(R_i) = B(R)/k$, for all i
- Goal: each R_i should fit in main memory:
 $B(R_i) \leq M$

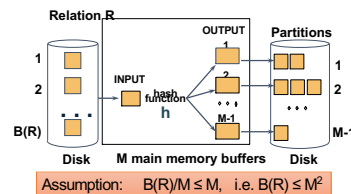
How do we choose k ?

January 27, 2020 CSE 444 - Winter 2020 74

74

Partitioned Hash Algorithms

- We choose $k = M-1$ Each bucket has size approx.
 $B(R)/(M-1) \approx B(R)/M$

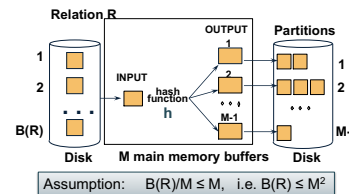


January 27, 2020 CSE 444 - Winter 2020 75

75

Partitioned Hash Algorithms

- We choose $k = M-1$ Each bucket has size approx.
 $B(R)/(M-1) \approx B(R)/M$



CSE 444 - Winter 2019

57

76

Grace-Join

 $R \bowtie S$

Note: grace-join is
also called
partitioned hash-join

January 27, 2020 CSE 444 - Winter 2020 77

77

Grace-Join

R \bowtie S

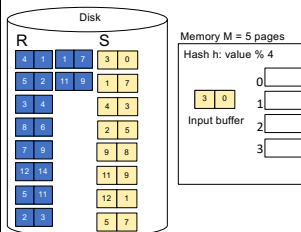
- Step 1:
 - Hash S into M-1 buckets
 - Send all buckets to disk
- Step 2
 - Hash R into M-1 buckets
 - Send all buckets to disk
- Step 3
 - Join every pair of buckets

Note: grace-join is
also called
partitioned hash-join

January 27, 2020 CSE 444 - Winter 2020 78

Partitioned Hash-Join Example

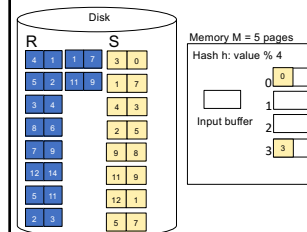
Step 1: Read relation S one page at a time and hash into M-1 (=4 buckets)



January 27, 2020 CSE 444 - Winter 2020 79

Partitioned Hash-Join Example

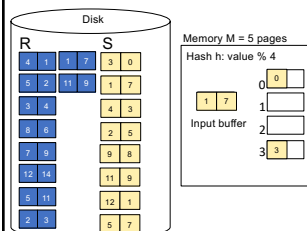
Step 1: Read relation S one page at a time and hash into the 4 buckets



January 27, 2020 CSE 444 - Winter 2020 80

Partitioned Hash-Join Example

Step 1: Read relation S one page at a time and hash into the 4 buckets

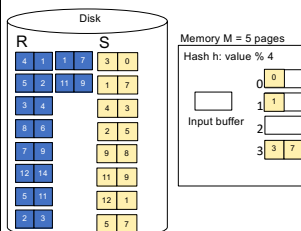


January 27, 2020 CSE 444 - Winter 2020 81

78

Partitioned Hash-Join Example

Step 1: Read relation S one page at a time and hash into the 4 buckets

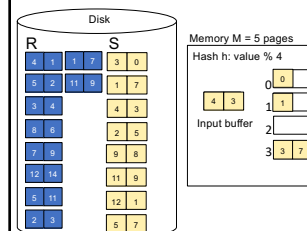


January 27, 2020 CSE 444 - Winter 2020 81

79

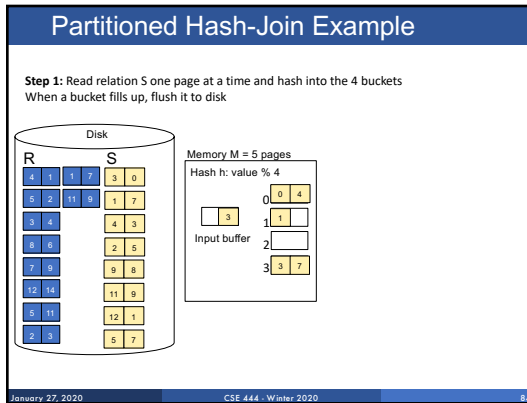
Partitioned Hash-Join Example

Step 1: Read relation S one page at a time and hash into the 4 buckets

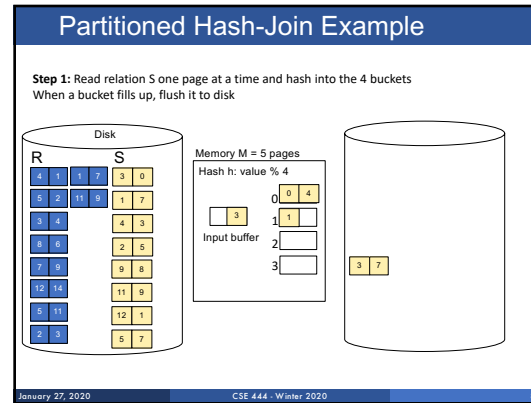


January 27, 2020 CSE 444 - Winter 2020 81

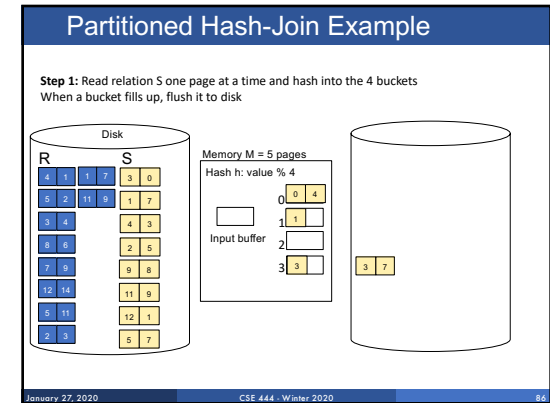
80



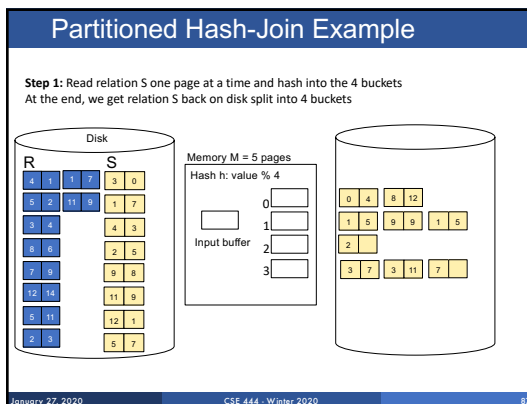
84



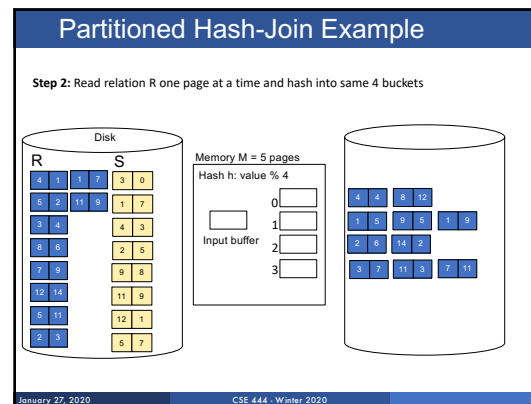
85



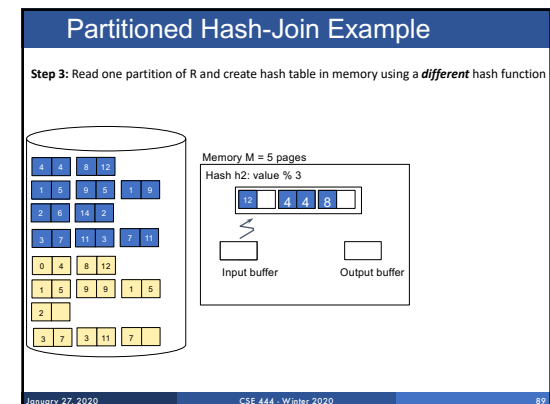
86



87



88



89

Partitioned Hash-Join Example

Step 4: Scan matching partition of S and probe the hash table
 Step 5: Repeat for all the buckets
 Total cost: $3B(R) + 3B(S)$

Memory $M = 5$ pages
 Hash $h2$: value % 3

Input buffer Output buffer

January 27, 2020 CSE 444 - Winter 2020 90

90

Grace-Join

Original Relation INPUT OUTPUT Partitions

Disk B main memory buffers Disk 1 2 ... M-1

Partition both relations using hash fn h . R tuples in partition i will only match S tuples in partition i .

January 27, 2020 CSE 444 - Winter 2020 91

91

Grace-Join

Original Relation INPUT OUTPUT Partitions

Disk B main memory buffers Disk 1 2 ... M-1

Partition both relations using hash fn h . R tuples in partition i will only match S tuples in partition i .

Read in a partition of R, hash it using $h2 (<= h1)$. Scan matching partition of S, search for matches.

Partitions of R & S Hash table for partition $S_i (< M-1$ pages) Join Result

Disk B main memory buffers Disk

January 27, 2020 CSE 444 - Winter 2020 92

92

Grace Join

- Cost: $3B(R) + 3B(S)$
- Assumption: $\min(B(R), B(S)) \leq M^2$

January 27, 2020 CSE 444 - Winter 2020 93

93

Hybrid Hash Join Algorithm

- Partition S into k buckets
 - t buckets S_1, \dots, S_t stay in memory
 - $k-t$ buckets S_{t+1}, \dots, S_k to disk
- Partition R into k buckets
 - First t buckets join immediately with S
 - Rest $k-t$ buckets go to disk
- Finally, join $k-t$ pairs of buckets:
 - $(R_{t+1}, S_{t+1}), (R_{t+2}, S_{t+2}), \dots, (R_k, S_k)$

January 27, 2020 CSE 444 - Winter 2020 94

94

Summary of External Join Algorithms

- Block Nested Loop: $B(S) + B(R) \cdot B(S) / (M-1)$
- Index Join: $B(R) + T(R)B(S) / V(S, a)$ (unclustered)
- Partitioned Hash: $3B(R) + 3B(S)$
 - $\min(B(R), B(S)) \leq M^2$
- Merge Join: $3B(R) + 3B(S)$
 - $B(R) + B(S) \leq M^2$

January 27, 2020 CSE 444 - Winter 2020 109

109

Summary of Query Execution

- For each logical query plan
 - There exist many physical query plans
 - Each plan has a different cost
 - Cost depends on the data
- Additionally, for each query
 - There exist several logical plans
- Next lecture: query optimization
 - How to compute the cost of a complete plan?
 - How to pick a good query plan for a query?

January 27, 2020

CSE 444 - Winter 2020

110

110