

Database System Internals Join Algorithms (cont.)

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

January 29, 2021 CSE 444 - Winter 2020 1

1

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)
- Merge Join: $3B(R) + 3B(S)$
 - $B(R) + B(S) \leq M^2$
- Partitioned Hash Join: (coming up next)

January 29, 2021 CSE 444 - Winter 2020 3

3

Partitioned Hash Algorithms

- Partition R it into k buckets:
 $R_1, R_2, R_3, \dots, R_k$

January 29, 2021 CSE 444 - Winter 2020 4

4

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 29, 2021 CSE 444 - Winter 2020 5

5

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 29, 2021

CSE 444 - Winter 2020

6

6

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 29, 2021

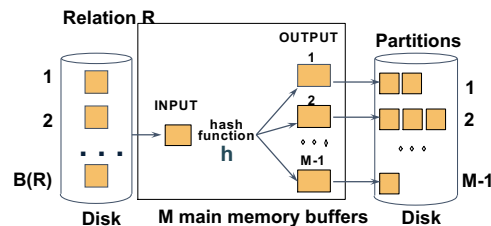
CSE 444 - Winter 2020

7

7

Partitioned Hash Algorithms

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

Assumption: $B(R)/M \leq M$, i.e. $B(R) \leq M^2$

January 29, 2021

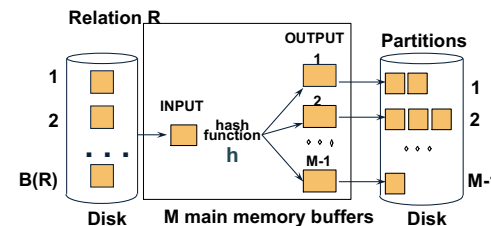
CSE 444 - Winter 2020

8

8

Partitioned Hash Algorithms

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

Assumption: $B(R)/M \leq M$, i.e. $B(R) \leq M^2$

CSE 444 - Winter 2019

57

9

Partitioned Hash Join (Grace-Join)

$R \bowtie S$

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

January 29, 2021

CSE 444 - Winter 2020

10

10

Partitioned Hash Join (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 29, 2021

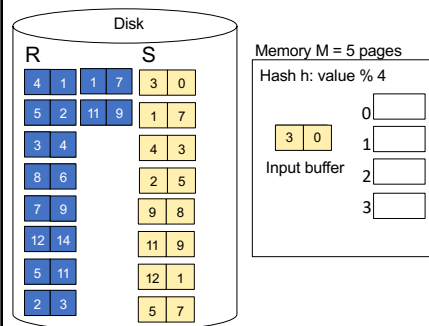
CSE 444 - Winter 2020

11

11

Partitioned Hash-Join Example

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



January 29, 2021

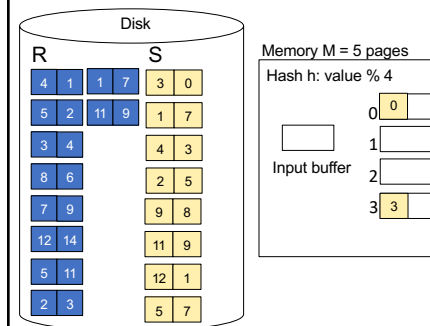
CSE 444 - Winter 2020

12

12

Partitioned Hash-Join Example

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



January 29, 2021

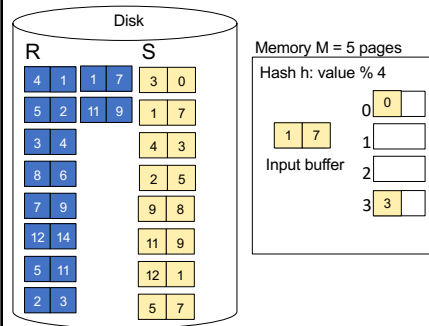
CSE 444 - Winter 2020

13

13

Partitioned Hash-Join Example

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



January 29, 2021

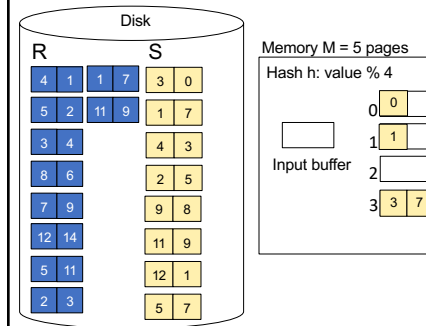
CSE 444 - Winter 2020

14

14

Partitioned Hash-Join Example

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



January 29, 2021

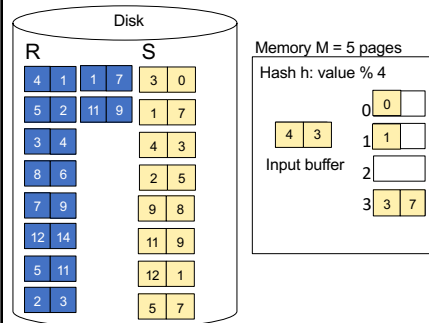
CSE 444 - Winter 2020

15

15

Partitioned Hash-Join Example

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



January 29, 2021

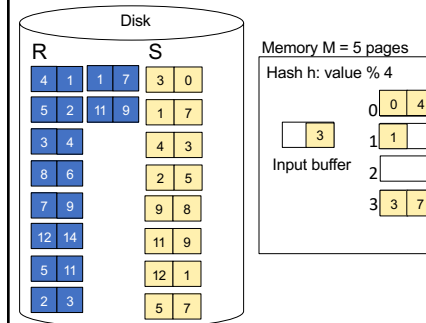
CSE 444 - Winter 2020

16

16

Partitioned Hash-Join Example

Step 1: Read relation S one page at a time and hash into the 4 buckets
When a bucket fills up, flush it to disk



January 29, 2021

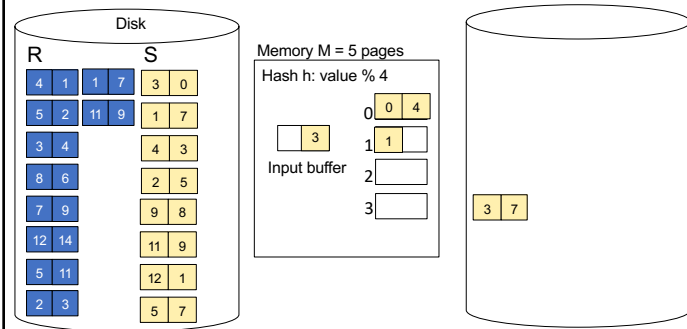
CSE 444 - Winter 2020

17

17

Partitioned Hash-Join Example

Step 1: Read relation S one page at a time and hash into the 4 buckets
When a bucket fills up, flush it to disk



January 29, 2021

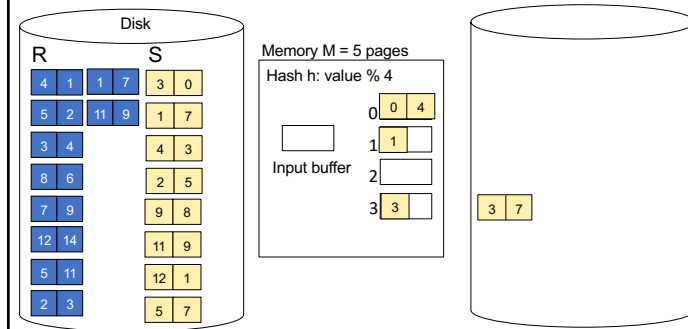
CSE 444 - Winter 2020

18

18

Partitioned Hash-Join Example

Step 1: Read relation S one page at a time and hash into the 4 buckets
When a bucket fills up, flush it to disk



January 29, 2021

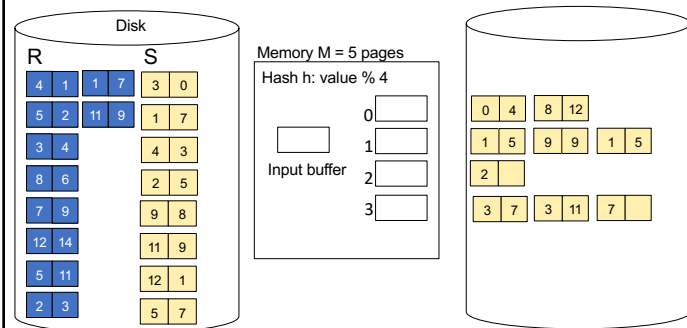
CSE 444 - Winter 2020

19

19

Partitioned Hash-Join Example

Step 1: Read relation S one page at a time and hash into the 4 buckets
At the end, we get relation S back on disk split into 4 buckets



January 29, 2021

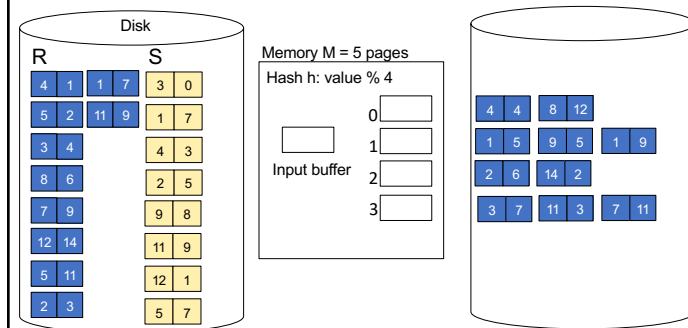
CSE 444 - Winter 2020

20

20

Partitioned Hash-Join Example

Step 2: Read relation R one page at a time and hash into same 4 buckets



January 29, 2021

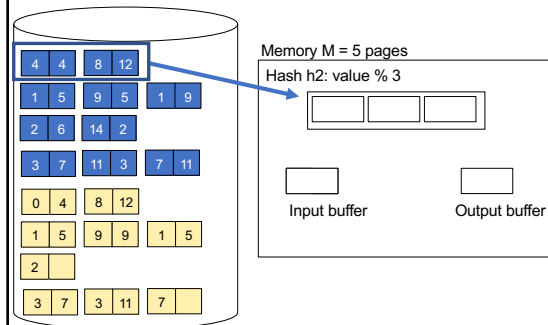
CSE 444 - Winter 2020

21

21

Partitioned Hash-Join Example

Step 3: Read one partition of R and create hash table in memory using a *different* hash function



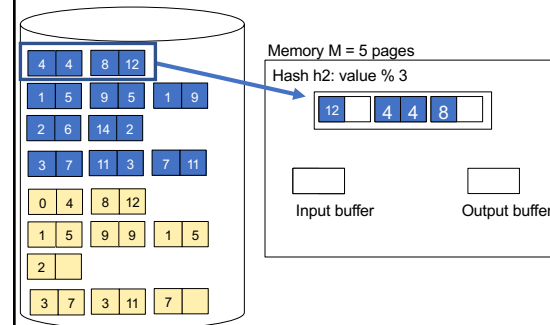
January 29, 2021

CSE 444 - Winter 2020

22

Partitioned Hash-Join Example

Step 3: Read one partition of R and create hash table in memory using a *different* hash function



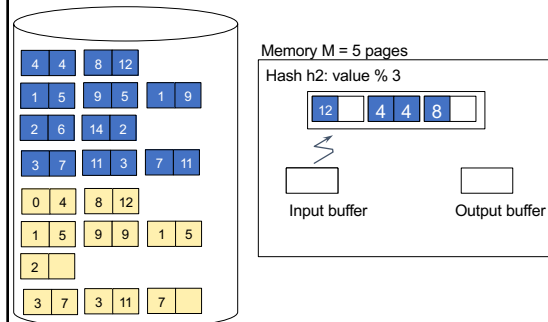
January 29, 2021

CSE 444 - Winter 2020

23

Partitioned Hash-Join Example

Step 3: Read one partition of R and create hash table in memory using a *different* hash function



January 29, 2021

CSE 444 - Winter 2020

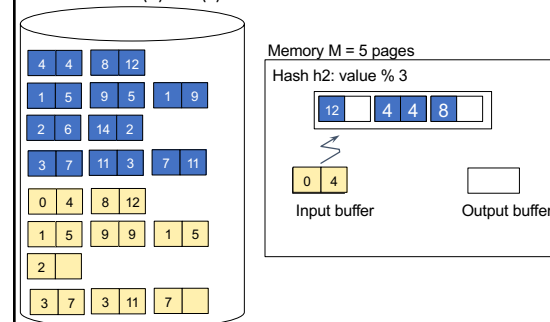
24

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: $3B(R) + 3B(S)$



January 29, 2021

CSE 444 - Winter 2020

25

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 h_2 : value % 3

Input buffer: 0 4
 Output buffer: 4 4

January 29, 2021 CSE 444 - Winter 2020 26

26

Partitioned Hash-Join

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

Original Relation INPUT hash function h OUTPUT 1 2 ... M-1 Partitions 1 2 ... M-1

Disk B main memory buffers Disk

January 29, 2021 CSE 444 - Winter 2020 27

27

Partitioned Hash-Join

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

Original Relation INPUT hash function h OUTPUT 1 2 ... M-1 Partitions 1 2 ... M-1

Disk B main memory buffers Disk

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

hash fn h_2 Input buffer for R_i Output buffer

Disk B main memory buffers Disk

❖ Read in a partition of R, hash it using h_2 ($\neq h_1$). Scan matching partition of S, search for matches.

January 29, 2021 CSE 444 - Winter 2020 28

28

Partitioned Hash-Join

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

January 29, 2021 CSE 444 - Winter 2020 29

29

Hybrid Hash Join Algorithm (see book)

- 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 29, 2021

CSE 444 - Winter 2020

30

30

Before We Go Into Query Plan Costs... How do Updates Work? (Insert/Delete)

January 29, 2021

CSE 444 - Winter 2020

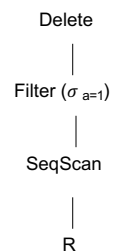
45

45

Example Using Delete

delete from R where a=1;

Query plan



In SimpleDB, the Delete Operator calls `BufferPool.deleteTuple()`

Why not call `HeapFile.deleteTuple()` directly?

Because there could also be indexes. Need some entity that will decide all the structures from where tuple needs to be deleted

BufferPool then calls `HeapFile.deleteTuple()`

January 29, 2021

CSE 444 - Winter 2020

46

46

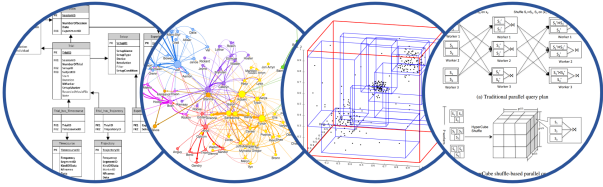
Pushing Updates to Disk

- When **inserting a tuple**, HeapFile inserts it on a page but does not write the page to disk
- When **deleting a tuple**, HeapFile deletes tuple from a page but does not write the page to disk
- The buffer manager worries when to write pages to disk (and when to read them from disk)
- When need to **add new page** to file, HeapFile adds page to file on disk and then reads it through buffer manager

January 29, 2021

47

47



Database System Internals Query Plan Costs

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

January 29, 2021 CSE 444 - Winter 2020 48

48

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 29, 2021 CSE 444 - Winter 2020 49

49

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 29, 2021 CSE 444 - Winter 2020 50

50

A Note About Skew

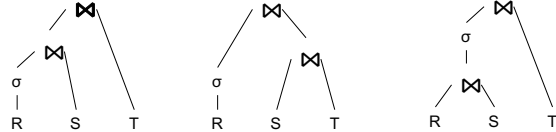
- Previously shown 2 pass join algorithms do not work for heavily skewed data
- For a sort-merge join, the maximum number of tuples with a particular join attribute should be the number of tuples per page:
 - This often isn't the case: would need multiple passes

January 29, 2021 CSE 444 - Winter 2020 52

52

Query Optimization Summary

Goal: find a physical plan that has minimal cost



What is the cost of a plan?

For each operator, cost is function of CPU, IO, network bw

$$\text{Total_Cost} = \text{CPUCost} + w_{\text{IO}} \text{IOCost} + w_{\text{BW}} \text{BWCost}$$

Cost of plan is total for all operators

In this class, we look only at IO

January 29, 2021

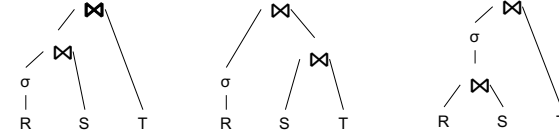
CSE 444 - Winter 2020

53

53

Query Optimization Summary

Goal: find a physical plan that has minimal cost



January 29, 2021

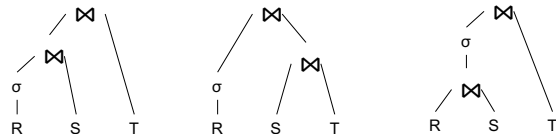
CSE 444 - Winter 2020

54

54

Query Optimization Summary

Goal: find a physical plan that has minimal cost



Know how to compute cost if know cardinalities

January 29, 2021

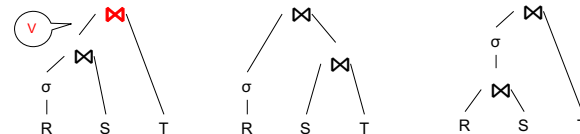
CSE 444 - Winter 2020

55

55

Query Optimization Summary

Goal: find a physical plan that has minimal cost



Know how to compute cost if know cardinalities

January 29, 2021

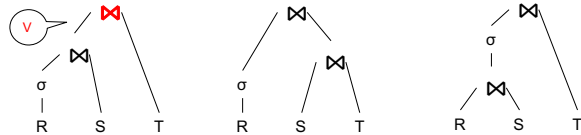
CSE 444 - Winter 2020

56

56

Query Optimization Summary

Goal: find a physical plan that has minimal cost



Know how to compute cost if know cardinalities

- Eg. $\text{Cost}(V \bowtie T) = 3B(V) + 3B(T)$
- $B(V) = T(V) / \text{PageSize}$
- $T(V) = T(\sigma(R) \bowtie S)$

January 29, 2021

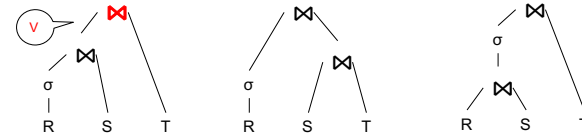
CSE 444 - Winter 2020

57

57

Query Optimization Summary

Goal: find a physical plan that has minimal cost



Know how to compute cost if know cardinalities

- Eg. $\text{Cost}(V \bowtie T) = 3B(V) + 3B(T)$
- $B(V) = T(V) / \text{PageSize}$
- $T(V) = T(\sigma(R) \bowtie S)$

Cardinality estimation problem: e.g. estimate $T(\sigma(R) \bowtie S)$

January 29, 2021

CSE 444 - Winter 2020

58

58

Database Statistics

- **Collect** statistical summaries of stored data
- **Estimate size** (=cardinality) in a bottom-up fashion
 - This is the most difficult part, and still inadequate in today's query optimizers
- **Estimate cost** by using the estimated size
 - Hand-written formulas, similar to those we used for computing the cost of each physical operator

January 29, 2021

CSE 444 - Winter 2020

59

59

Database Statistics

- Number of tuples (cardinality) $T(R)$
- Indexes, number of keys in the index $V(R,a)$
- Number of physical pages $B(R)$
- Statistical information on attributes
 - Min value, Max value, $V(R,a)$
- Histograms
- **Collection approach: periodic, using sampling**

January 29, 2021

CSE 444 - Winter 2020

60

60

Size Estimation Problem

```
Q = SELECT list
    FROM R1, ..., Rn
    WHERE cond1 AND cond2 AND ... AND condk
```

Given $T(R_1), T(R_2), \dots, T(R_n)$
Estimate $T(Q)$

How can we do this ? Note: doesn't have to be exact.

January 29, 2021

CSE 444 - Winter 2020

61

61

Size Estimation Problem

```
Q = SELECT list
    FROM R1, ..., Rn
    WHERE cond1 AND cond2 AND ... AND condk
```

Remark: $T(Q) \leq T(R_1) \times T(R_2) \times \dots \times T(R_n)$

January 29, 2021

CSE 444 - Winter 2020

62

62

Size Estimation Problem

```
Q = SELECT list
    FROM R1, ..., Rn
    WHERE cond1 AND cond2 AND ... AND condk
```

Remark: $T(Q) \leq T(R_1) \times T(R_2) \times \dots \times T(R_n)$

Key idea: each condition reduces the size of $T(Q)$ by some factor, called **selectivity factor**

January 29, 2021

CSE 444 - Winter 2020

63

63

Selectivity Factor

- Each condition **cond** reduces the size by some factor called **selectivity factor**
- Assuming independence, **multiply** the selectivity factors

January 29, 2021

CSE 444 - Winter 2020

64

64

Example

R(A,B)
S(B,C)
T(C,D)

```
Q = SELECT *
    FROM R, S, T
    WHERE R.B=S.B and S.C=T.C and R.A<40
```

$T(R) = 30k$, $T(S) = 200k$, $T(T) = 10k$

Selectivity of $R.B = S.B$ is $1/3$
Selectivity of $S.C = T.C$ is $1/10$
Selectivity of $R.A < 40$ is $1/2$

Q: What is the estimated size of the query output $T(Q)$?

January 29, 2021

CSE 444 - Winter 2020

65

65

Example

R(A,B)
S(B,C)
T(C,D)

```
Q = SELECT *
    FROM R, S, T
    WHERE R.B=S.B and S.C=T.C and R.A<40
```

$T(R) = 30k$, $T(S) = 200k$, $T(T) = 10k$

Selectivity of $R.B = S.B$ is $1/3$
Selectivity of $S.C = T.C$ is $1/10$
Selectivity of $R.A < 40$ is $1/2$

Q: What is the estimated size of the query output $T(Q)$?

A: $T(Q) = 30k * 200k * 10k * 1/3 * 1/10 * 1/2 = 10^{12}$

January 29, 2021

CSE 444 - Winter 2020

66

66

Selectivity Factors for Conditions

- $A = c$ $/ * \sigma_{A=c}(R) */$
 - Selectivity = $1/V(R,A)$

January 29, 2021

CSE 444 - Winter 2020

67

67

Selectivity Factors for Conditions

- $A = c$ $/ * \sigma_{A=c}(R) */$
 - Selectivity = $1/V(R,A)$
- $A < c$ $/ * \sigma_{A< c}(R) */$
 - Selectivity = $(c - \text{Low}(R, A)) / (\text{High}(R, A) - \text{Low}(R, A))$

January 29, 2021

CSE 444 - Winter 2020

68

68

Selectivity Factors for Conditions

- $A = c$ /* $\sigma_{A=c}(R)$ */
 - Selectivity = $1/V(R,A)$
- $A < c$ /* $\sigma_{A < c}(R)$ */
 - Selectivity = $(c - \text{Low}(R,A)) / (\text{High}(R,A) - \text{Low}(R,A))$
- $A = B$ /* $R \bowtie_{A=B} S$ */
 - Selectivity = $1 / \max(V(R,A), V(S,A))$
 - (will explain next)

January 29, 2021

CSE 444 - Winter 2020

69

69

Assumptions

- Containment of values: if $V(R,A) \leq V(S,B)$, then all values $R.A$ occur in $S.B$
 - Note: this indeed holds when A is a foreign key in R , and B is a key in S
- Preservation of values: for any other attribute C , $V(R \bowtie_{A=B} S, C) = V(R, C)$ (or $V(S, C)$)
 - Note: we don't need this to estimate the size of the join, but we need it in estimating the next operator

January 29, 2021

CSE 444 - Winter 2020

70

70

Selectivity of $R \bowtie_{A=B} S$

Assume $V(R,A) \leq V(S,B)$

- A tuple t in R joins with $T(S)/V(S,B)$ tuple(s) in S
- Hence $T(R \bowtie_{A=B} S) = T(R) T(S) / V(S,B)$

$$T(R \bowtie_{A=B} S) = T(R) T(S) / \max(V(R,A), V(S,B))$$

January 29, 2021

CSE 444 - Winter 2020

71

71

Complete Example

Supplier(sno, sname, scity, sstate)
Supply(sno, pno, quantity)

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sno = y.sno
and y.pno = 2
and x.scity = 'Seattle'
and x.sstate = 'WA'
```

- Some statistics
 - $T(\text{Supplier}) = 1000$ records
 - $T(\text{Supply}) = 10,000$ records
 - $B(\text{Supplier}) = 100$ pages
 - $B(\text{Supply}) = 100$ pages
 - $V(\text{Supplier}, \text{scity}) = 20$, $V(\text{Suppliers}, \text{state}) = 10$
 - $V(\text{Supply}, \text{pno}) = 2,500$
 - Both relations are clustered
- $M = 11$

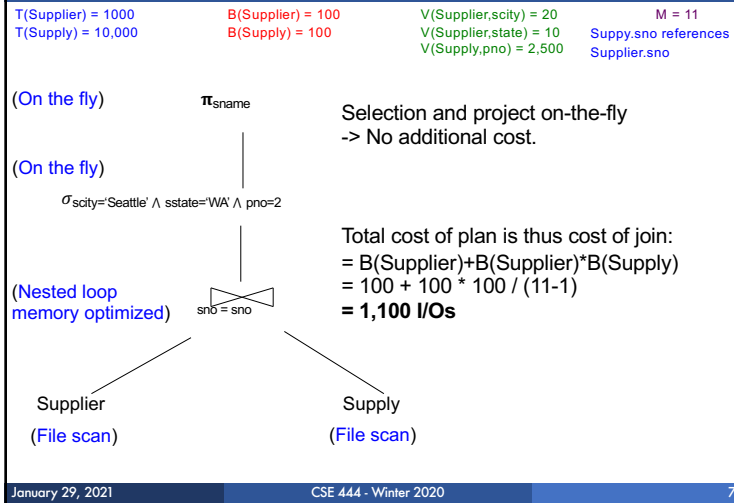
January 29, 2021

CSE 444 - Winter 2020

73

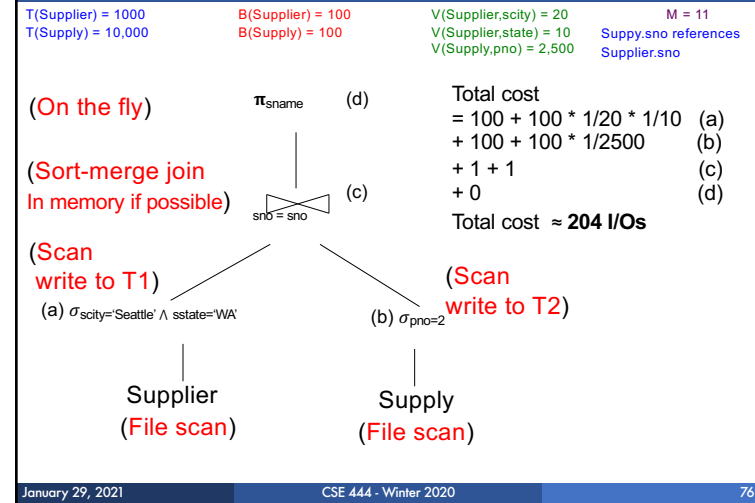
73

Physical Query Plan 1



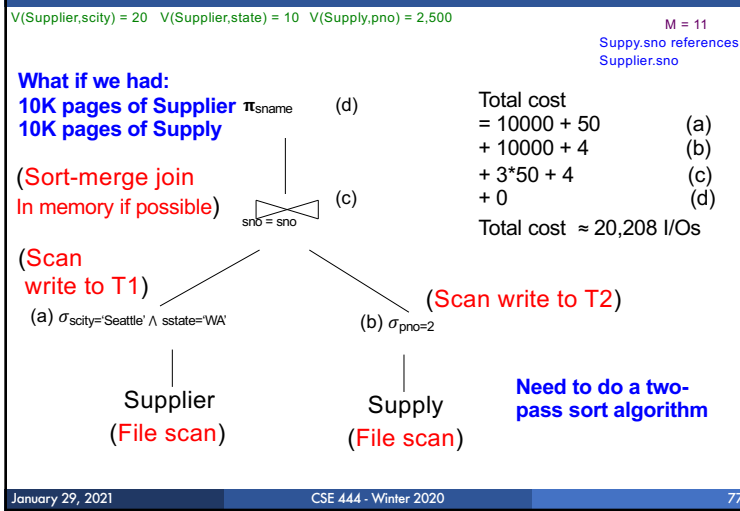
75

Physical Query Plan 2



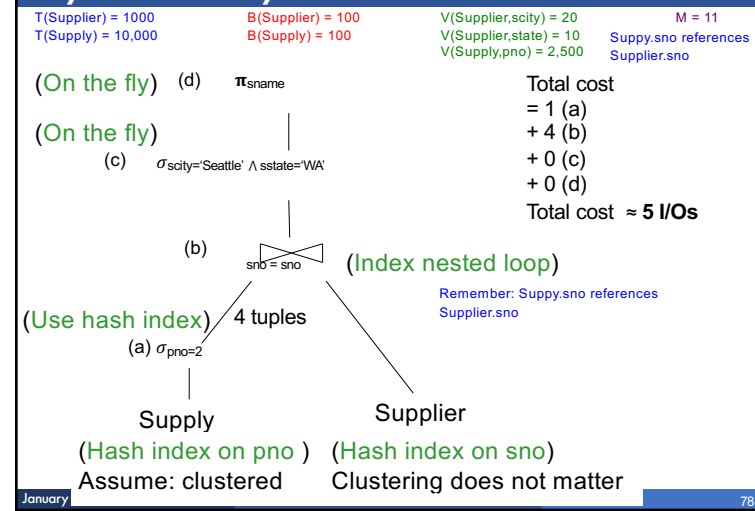
76

Plan 2 with Different Numbers



77

Physical Query Plan 3



78

Histograms

- Statistics on data maintained by the RDBMS
- Makes size estimation much more accurate (hence, cost estimations are more accurate)

January 29, 2021

CSE 444 - Winter 2020

79

79

Histograms

Employee(ssn, name, age)

$T(\text{Employee}) = 25000$, $V(\text{Employee, age}) = 50$
 $\min(\text{age}) = 19$, $\max(\text{age}) = 68$

$\sigma_{\text{age}=48}(\text{Employee}) = ?$ $\sigma_{\text{age}>28 \text{ and age}<35}(\text{Employee}) = ?$

January 29, 2021

CSE 444 - Winter 2020

80

80

Histograms

Employee(ssn, name, age)

$T(\text{Employee}) = 25000$, $V(\text{Employee, age}) = 50$
 $\min(\text{age}) = 19$, $\max(\text{age}) = 68$

$\sigma_{\text{age}=48}(\text{Employee}) = ?$ $\sigma_{\text{age}>28 \text{ and age}<35}(\text{Employee}) = ?$



Estimate = $25000 / 50 = 500$



Estimate = $25000 * 6 / 50 = 3000$

January 29, 2021

CSE 444 - Winter 2020

81

81

Histograms

Employee(ssn, name, age)

$T(\text{Employee}) = 25000$, $V(\text{Employee, age}) = 50$
 $\min(\text{age}) = 19$, $\max(\text{age}) = 68$

$\sigma_{\text{age}=48}(\text{Employee}) = ?$ $\sigma_{\text{age}>28 \text{ and age}<35}(\text{Employee}) = ?$

Age:	0-20	20-29	30-39	40-49	50-59	> 60
Tuples	200	800	5000	12000	6500	500

January 29, 2021

CSE 444 - Winter 2020

82

82

Histograms

Employee(ssn, name, age)

$T(\text{Employee}) = 25000$, $V(\text{Employee}, \text{age}) = 50$
 $\min(\text{age}) = 19$, $\max(\text{age}) = 68$

$\sigma_{\text{age}=48}(\text{Employee}) = ?$ $\sigma_{\text{age}>28 \text{ and } \text{age}<35}(\text{Employee}) = ?$

Age:	0..20	20..29	30-39	40-49	50-59	> 60
Tuples	200	800	5000	12000	6500	500

Estimate = 1200

Estimate = $1 \cdot 80 + 5 \cdot 500 = 2580$

January 29, 2021

CSE 444 - Winter 2020

83

83

Types of Histograms

- How should we determine the bucket boundaries in a histogram?

January 29, 2021

CSE 444 - Winter 2020

84

84

Types of Histograms

- How should we determine the bucket boundaries in a histogram?

- Eq-Width
- Eq-Depth
- Compressed
- V-Optimal histograms

January 29, 2021

CSE 444 - Winter 2020

85

85

Histograms

Employee(ssn, name, age)

Eq-width:

Age:	0..20	20..29	30-39	40-49	50-59	> 60
Tuples	200	800	5000	12000	6500	500

Eq-depth:

Age:	0-33	33-38	38-43	43-45	45-54	> 54
Tuples	1800	2000	2100	2200	1900	1800

Compressed: store separately highly frequent values: (48,1900)

January 29, 2021

CSE 444 - Winter 2020

86

86

V-Optimal Histograms

- Defines bucket boundaries in an optimal way, to minimize the error over all point queries
- Computed rather expensively, using dynamic programming
- Modern databases systems use V-optimal histograms or some variations

January 29, 2021

CSE 444 - Winter 2020

87

87

Difficult Questions on Histograms

- Small number of buckets
 - Hundreds, or thousands, but not more
 - WHY ?
- Not updated during database update, but recomputed periodically
 - WHY ?
- Multidimensional histograms rarely used
 - WHY ?

January 29, 2021

CSE 444 - Winter 2020

88

88

Difficult Questions on Histograms

- Small number of buckets
 - Hundreds, or thousands, but not more
 - WHY? All histograms are kept in main memory during query optimization; plus need fast access
- Not updated during database update, but recomputed periodically
 - WHY? Histogram update creates a write conflict; would dramatically slow down transaction throughput
- Multidimensional histograms rarely used
 - WHY? Too many possible multidimensional histograms, unclear which ones to choose

January 29, 2021

CSE 444 - Winter 2020

89

89