

Name

CWID

Quiz 2

Due April 26th, 2017

CS525 - Advanced Database Organization Solutions

Please leave this empty!

1

2

3

4

5

6

7

Sum

Instructions

- Multiple choice questions are graded in the following way: You get points for correct answers and points subtracted for wrong answers. The minimum points for each questions is **0**. For example, assume there is a multiple choice question with 6 answers - each may be correct or incorrect - and each answer gives 1 point. If you answer 3 questions correct and 3 incorrect you get 0 points. If you answer 4 questions correct and 2 incorrect you get 2 points. . . .
- For your convenience the number of points for each part and questions are shown in parenthesis.
- There are 3 parts in this quiz
 1. Disk Organization and Buffering
 2. Index Structures
 3. Result Size Estimations
 4. I/O Cost Estimation
 5. Schedules
 6. ARIES (Optional)
 7. Physical Optimization (Optional)

Part 1 Disk Organization and Buffering (Total: 15 Points)

Question 1.1 Page Replacement Clock (15 Points)

Consider a buffer pool with 3 pages using the **Clock** page replacement strategy. Initially the buffer pool is in the state shown below. We use the following notation $flag[page]_{fix}^{dirty}$ to denote the state of each buffer frame. $page$ is the number of the page in the frame, fix is its fix count, $dirty$ is indicating with an Asterix that the page is dirty, and $flag$ is the reference bit used by the Clock algorithm. E.g., $^1[5]_2^*$ denotes that the frame stores page 5 with a fix count 2, that the page is dirty, and that the reference bit is set to 1. Recall that Clock uses a pointer S that points to the current page frame (the one to be checked for replacement next). The page frame S is pointing to is indicated by \downarrow . In your solution draw an arrow to the page frame that S is pointing to.

Current Buffer State

$$^1[3]_0^* \quad \downarrow \quad ^1[10]_0 \quad ^0[6]_0 \quad ^1[4]_1$$

Execute the following requests and write down state of the buffer pool after each request.

- p stands for pin
- u for unpin
- d for marking a page as dirty

$$p(4), d(4), p(1), p(2), u(2), p(9), p(7)$$

Solution

p(4)

$$^1[3]_0^* \quad \downarrow \quad ^1[10]_0 \quad ^0[6]_0 \quad ^1[4]_2$$

d(4)

$$^1[3]_0^* \quad \downarrow \quad ^1[10]_0 \quad ^0[6]_0 \quad ^1[4]_2^*$$

p(1)

$$^1[3]_0^* \quad ^0[10]_0 \quad ^1[1]_1 \quad \downarrow \quad ^1[4]_2^*$$

p(2)

$$^0[3]_0^* \quad ^1[2]_1 \quad \downarrow \quad ^1[1]_1 \quad ^0[4]_2^*$$

u(2)

$$^0[3]_0^* \quad ^1[2]_0 \quad \downarrow \quad ^1[1]_1 \quad ^0[4]_2^*$$

p(9)

$$^1[9]_1 \quad \downarrow \quad ^1[2]_0 \quad ^0[1]_1 \quad ^0[4]_2^*$$

p(7)

$$^0[9]_1 \quad ^1[7]_1 \quad \downarrow \quad ^0[1]_1 \quad ^0[4]_2^*$$

Part 2 Index Structures (Total: 25 Points)

Assume that you have the following table:

Item		
name	age	salary
Peter Petersen	35	10,000
Gert Gertsen	20	40,000
Heinz Heinzen	19	100,000
Gertrud Gertsen	19	45,000
Astrid Lundgren	19	110,000
Pferdegert	20	50,000
Heinz Bert	65	38,000

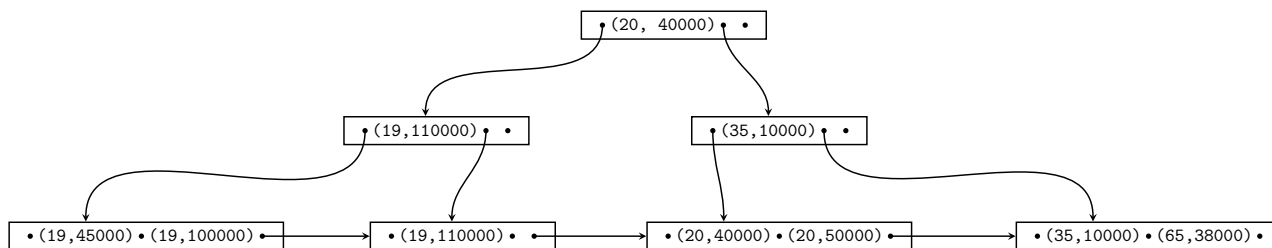
Question 2.1 Construction (9 Points)

Create a B+-tree for table *Item* over attributes *age* and *salary* with $n = 2$ (up to two keys per node). You should start with an empty B+-tree and insert the keys in the order shown in the table above. Write down the resulting B+-tree after each step.

When splitting or merging nodes follow these conventions:

- **Leaf Split:** In case a leaf node needs to be split during insertion and n is even, the left node should get the extra key. E.g, if $n = 2$ and we insert a key 4 into a node $[1,5]$, then the resulting nodes should be $[1,4]$ and $[5]$. For odd values of n we can always evenly split the keys between the two nodes. In both cases the value inserted into the parent is the smallest value of the right node.
- **Non-Leaf Split:** In case a non-leaf node needs to be split and n is odd, we cannot split the node evenly (one of the new nodes will have one more key). In this case the “middle” value inserted into the parent should be taken from the right node. E.g., if $n = 3$ and we have to split a non-leaf node $[1,3,4,5]$, the resulting nodes would be $[1,3]$ and $[5]$. The value inserted into the parent would be 4.
- **Node Underflow:** In case of a node underflow you should first try to redistribute values from a sibling and only if this fails merge the node with one of its siblings. Both approaches should prefer the left sibling. E.g., if we can borrow values from both the left and right sibling, you should borrow from the left one.

Solution

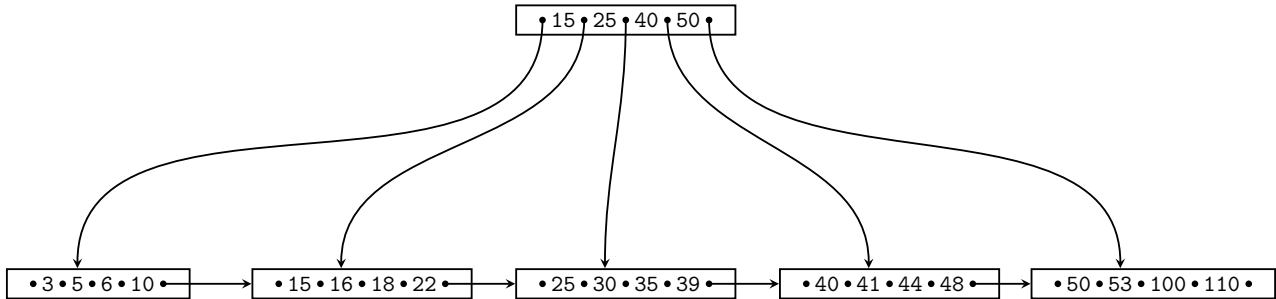


Question 2.2 Operations (9 Points)

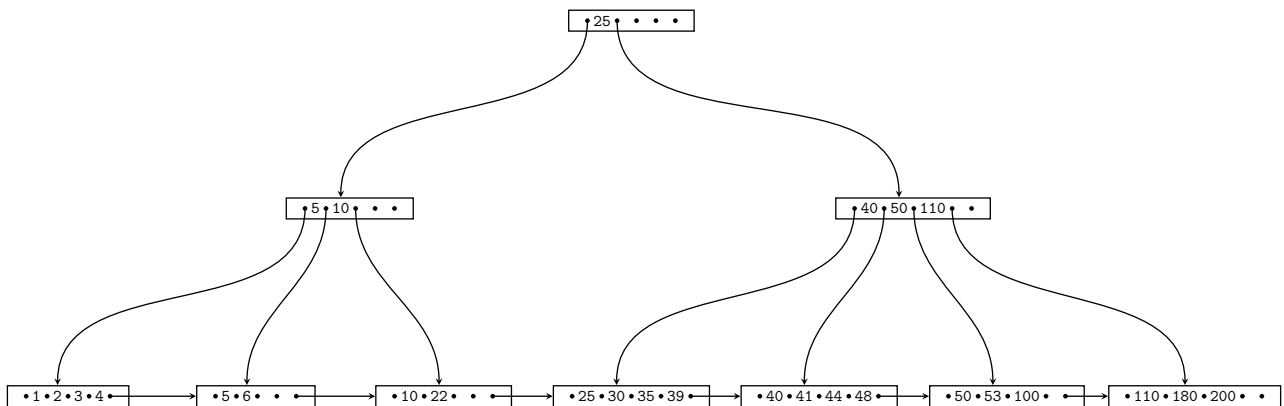
Given is the B+-tree shown below ($n = 4$). Execute the following operations and write down the resulting B+-tree after each operation:

delete(15), delete(16), delete(18), insert(1), insert(2), insert(4), insert(200), insert(180)

Use the conventions for splitting and merging introduced in the previous question.



Solution



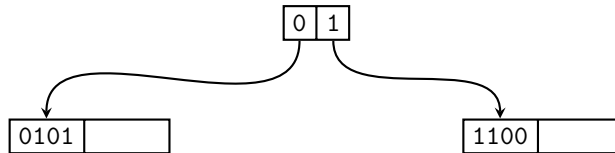
Question 2.3 Extensible Hashing (7 Points)

Consider the extensible Hash index shown below that is the result of inserting values 6 and 4. Each page holds two keys. Execute the following operations

`insert(1), insert(8), insert(3), insert(2), insert(7)`

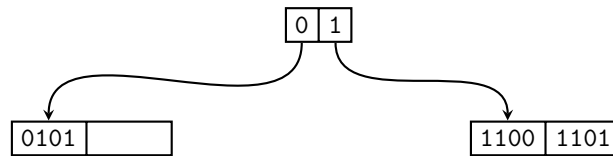
and write down the resulting index after each operation. Assume the hash function is defined as:

x	h(x)
0	1000
1	1101
2	0111
3	0000
4	1100
5	0100
6	0101
7	1001
8	1110

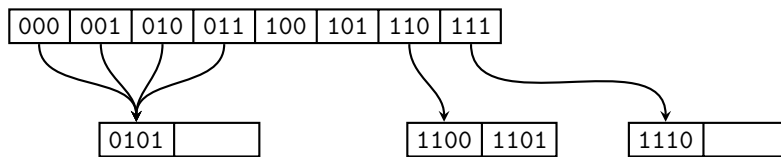


Solution

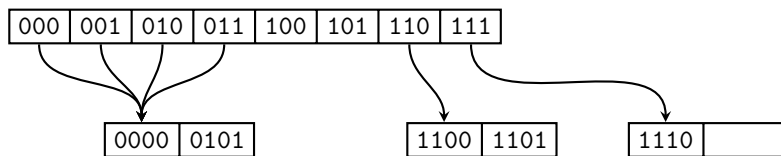
`insert(1)`



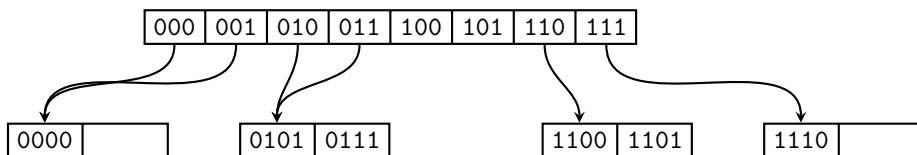
`insert(8)`



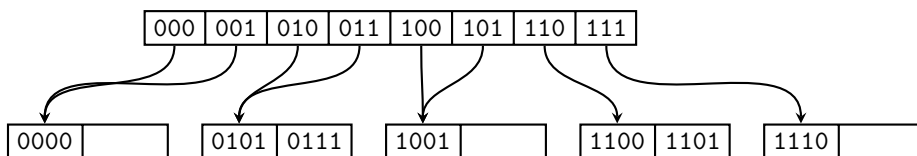
`insert(3)`



`insert(2)`



`insert(7)`



Part 3 Result Size Estimations (Total: 20 Points)

Consider a table *book* with attributes ISBN, title, author, edition (primary key is ISBN), a table *library* with loc, budget, public (primary key is loc), and a table *catalog* with attributes library and book. *catalog.library* is a foreign key to relation *library*. Attribute *book* of relation *catalog* is a foreign key to relation *book*. Given are the following statistics:

$$\begin{array}{lll} T(\textit{book}) = 100,000 & T(\textit{library}) = 100 & T(\textit{catalog}) = 200,000 \\ V(\textit{book}, \textit{ISBN}) = 100,000 & V(\textit{library}, \textit{loc}) = 100 & V(\textit{catalog}, \textit{library}) = 100 \\ V(\textit{book}, \textit{title}) = 50,000 & V(\textit{library}, \textit{budget}) = 40 & V(\textit{catalog}, \textit{book}) = 90,000 \\ V(\textit{book}, \textit{author}) = 30,000 & V(\textit{library}, \textit{public}) = 2 & \\ V(\textit{book}, \textit{edition}) = 15 & & \end{array}$$

Question 3.1 Estimate Result Size (4 Points)

Estimate the number of result tuples for the query $q = \sigma_{\textit{public}=\textit{true}}(\textit{library})$ using the first assumption presented in class (values used in queries are uniformly distributed within the active domain).

Solution

$$T(q) = \frac{T(\textit{library})}{V(\textit{library}, \textit{public})} = \frac{100}{2} = 50$$

Question 3.2 Estimate Result Size (5 Points)

Estimate the number of result tuples for the query $q = \sigma_{\textit{title}=\textit{Faust} \wedge \textit{author}=\textit{Goethe}}(\textit{book})$ using the first assumption presented in class.

Solution

$$T(q) = \frac{T(\text{book})}{V(\text{book}, \text{title}) \cdot V(\text{book}, \text{author})} = \frac{100,000}{50,000 \cdot 30,000} = \frac{1}{15,000}$$

Question 3.3 Estimate Result Size (5 Points)

Estimate the number of result tuples for the query $q = \sigma_{\text{edition} \geq 2 \wedge \text{edition} \leq 4 \wedge \text{title} = \text{Databases}}(\text{book})$ using the first assumption presented in class. Assume that the minimal and maximal values in the `edition` attribute are 1 and 15,

Solution

$$T(q) = \frac{4 - 2 + 1}{\max(\text{book}, \text{edition}) - \min(\text{book}, \text{edition}) + 1} \cdot \frac{1}{50,000} \cdot T(\text{book}) = \frac{3 \cdot 100,000}{15 \cdot 50,000} = \frac{2}{5} = 0.4$$

Question 3.4 Estimate Result Size (6 Points)

Estimate the number of result tuples for the query

$q = \sigma_{\text{title} = \text{DatabaseIntroduction} \vee \text{title} = \text{DatabaseSystems}}(\text{book}) \bowtie_{\text{title} = \text{book}} \text{catalog} \bowtie_{\text{library} = \text{loc}} \sigma_{\text{budget} \leq 40}(\text{library})$

using the first assumption presented in class.

Solution

Let $q_1 = \sigma_{title=DatabaseIntroduction \vee title=DatabaseSystems}(book)$ and $q_2 = \sigma_{budget \leq 40}(library)$.
To estimate the selection result size q_1 :

$$\begin{aligned} T(q_1) &= (1 - [(1 - \frac{1}{V(book, title)}) \cdot (1 - \frac{1}{V(book, title)})]) \cdot T(book) \\ &= (1 - (1 - \frac{1}{50,000}) \cdot (1 - \frac{1}{50,000})) \cdot 10,000 \approx 0.4 \end{aligned}$$

To estimate the selection result size q_2 :

$$\begin{aligned} T(q_2) &= \frac{\max(library, budget) - 40 + 1}{\max(library, budget) - \min(library, budget) + 1} \cdot T(library) \\ &= \frac{70 - 40 + 1}{70 - 10 + 1} \cdot 100 = \frac{31}{61} \cdot 100 \approx 51 \end{aligned}$$

Now for the full query we get

$$\begin{aligned} T(q) &= \frac{T(q_1) \cdot T(catalog) \cdot T(q_2)}{\max(V(q_1, title), V(catalog, book)) \cdot \max(V(catalog, library), V(q_2, loc))} \\ &= \frac{0.4 \cdot 200,000 \cdot 51}{\max(2,90000) \cdot \max(100, 51)} = \frac{20,400}{9,000,000} \approx 0.45 \end{aligned}$$

Part 4 I/O Cost Estimation (Total: 20 Points)

Question 4.1 External Sorting (4 Points)

You have $M = 3$ memory pages available and should sort a relation R with $B(R) = 30,000,000$ blocks. Compute the number of I/Os necessary to sort R using the external merge sort algorithm introduced in class.

Solution

$$\begin{aligned} IO &= 2 \cdot B(R) \cdot (1 + \lceil \log_{M-1}(\frac{B(R)}{M}) \rceil) \\ &= 2 \cdot 30,000,000 \cdot (1 + 24) \\ &= 1,500,000,000 \end{aligned}$$

Question 4.2 External Sorting (4 Points)

You have $M = 201$ memory pages available and should sort a relation R with $B(R) = 201,000$ blocks. Compute the number of I/Os necessary to sort R using the external merge sort algorithm introduced in class.

Solution

$$\begin{aligned} IO &= 2 \cdot B(R) \cdot (1 + \lceil \log_{M-1}(\frac{B(R)}{M}) \rceil) \\ &= 2 \cdot 201,000 \cdot (1 + 2) \\ &= 1,206,000 \end{aligned}$$

Question 4.3 I/O Cost Estimation (6 = 2+2+2 Points)

Consider two relations R and S with $B(R) = 1,000,000$ and $B(S) = 2,000$. You have $M = 101$ memory pages available. Estimate the minimum number of I/O operations needed to join these two relations using **block-nested-loop join**, **merge-join** (the inputs are not sorted), and **hash-join**. You can assume that the hash function evenly distributes keys across buckets. Justify your result by showing the I/O cost estimation for each join method.

Solution

- **BNL**: S is smaller, thus, keep chunks of S in memory
 $\lceil \frac{B(S)}{M-1} \rceil \cdot [B(R) + \min(B(S), (M-1))] = 20 \cdot [1,000,000 + 100] = 20,002,000$ I/Os
- **MJ**: We can generate sorted runs of size 100. We need 2 merge passes for the sort for R and 1 merge passes for S . The last merge of R requires 99 pages, i.e., the number of sorted runs from R and S is not low enough to keep one page from each run of both R and S in memory. $7 \cdot B(R) + 3 \cdot B(S) = 7 \cdot 1,000,000 + 5 \cdot 2,000 = 7,010,000$ I/Os.
- **HJ**: We need 1 partitioning pass, because we can create 100 buckets and the bucket sizes of R and S will be 10,000 and 20 after one pass. The cost is $(2+1) \cdot (B(R) + B(S)) = 3 \cdot (1,000,000 + 2,000) = 3,006,000$ I/Os.

Question 4.4 I/O Cost Estimation (6 = 2+2+2 Points)

Consider two relations R and S with $B(R) = 2,000$ and $B(S) = 2,000$. You have $M = 21$ memory pages available. Compute the minimum number of I/O operations needed to join these two relations using **block-nested-loop join**, **merge-join** (the inputs are not sorted), and **hash-join**. You can assume that the hash function evenly distributes keys across buckets. Justify your result by showing the I/O cost estimation for each join method.

Solution

- **BNL**: R is smaller, thus, keep chunks of R in memory
 $\lceil \frac{B(R)}{M-1} \rceil \cdot [B(S) + \min(B(R), (M-1))] = 100 \cdot [2,000 + 100] = 210,000$ I/Os
- **MJ**: We can generate sorted runs of size 21 that means the number of sorted runs from R and S is low enough after 2 merge passes to keep one page from each run of both R and S in memory (5 runs from each). Thus, we need 2 merge passes for the sort, but can execute the last merge phase and join in one pass. $5 \cdot (B(R) + B(S)) = 5 \cdot (4,000) = 20,000$ I/Os.
- **HJ**: We need 2 partitioning passes, because we can create 20 buckets. The bucket sizes of R and S after the 2nd partitioning step will be 5. Thus, we can fit one bucket from R into memory to join with S (or vice versa). The cost is $(4+1) \cdot (B(R) + B(S)) = 5 \cdot (2,000 + 2,000) = 20,000$ I/Os.

Part 5 Schedules (Total: 20 Points)

Question 5.1 Schedule Classes (20 Points)

Indicate which of the following schedules belong to which class. Recall transaction operations are modelled as follows:

$w_1(A)$ transaction 1 wrote item A
 $r_1(A)$ transaction 1 read item A
 c_1 transaction 1 commits
 a_1 transaction 1 aborts

$S_1 = w_4(B), r_2(B), c_2, w_3(B), c_3, r_4(A), c_4$

$S_2 = w_2(A), r_4(B), r_3(E), w_2(C), c_2, r_1(E), w_4(A), w_4(B), c_4, r_3(D), w_3(C), w_3(D), c_3, r_1(D), c_1$

$S_3 = r_2(A), r_4(B), r_2(B), r_4(A), w_4(B), c_4, r_2(B), w_2(A), c_2$

$S_4 = r_3(B), w_3(B), r_3(A), w_4(A), r_2(B), w_2(B), c_2, r_4(B), c_4, r_3(A), c_3$

- S_1 is recoverable
- S_1 is cascade-less
- S_1 is strict
- S_1 is conflict-serializable
- S_1 is 2PL

- S_2 is recoverable
- S_2 is cascade-less
- S_2 is strict
- S_2 is conflict-serializable
- S_2 is 2PL

- S_3 is recoverable
- S_3 is cascade-less
- S_3 is strict
- S_3 is conflict-serializable
- S_3 is 2PL

- S_4 is recoverable
- S_4 is cascade-less
- S_4 is strict
- S_4 is conflict-serializable
- S_4 is 2PL

Part 6 Optional: ARIES (Total: 10 Optional Points)

Question 6.1 Recovery (10 Points)

Consider the state of the log and pages on disk shown below. For simplicity we do not show the actual undo/redo actions for updates, but instead show only the affected page. Assume a crash occurred after the last log entry. Answer the following questions:

1. **Analysis:** Write down the result of the analysis phase (RedoLSN, Transaction Table, Dirty Page Table)
2. **Redo:** Which pages will be loaded from disk during redo? Which pages will be modified during redo?
3. **Undo:** Write down the additional log entries that will be written during undo.

Log

LSN	Type	TID	PrevLSN	UndoNxtLSN	Data
1	begin	1	-	-	-
2	update	1	1	-	Page 4
3	begin	2	-	-	-
4	begin	3	-	-	-
5	update	2	4	-	Page 1
6	begin_cp	-	-	-	-
7	update	3	5	-	Page 1
8	update	3	7	-	Page 3
9	update	2	5	-	Page 2
10	commit	3	8	-	-
11	update	1	2	-	Page 2
12	update	2	9	-	Page 4
13	commit	2	12	-	-

Disk

PageID	PageLSN
1	7
2	9
3	0
4	2
5	0

Solution

(1):

RedoLSN: 2

Transaction Table: $\langle T_1, u, 11, - \rangle$

Dirty Page Table: $\langle 1, 5 \rangle, \langle 2, 9 \rangle, \langle 3, 8 \rangle, \langle 4, 2 \rangle$

(2):

All pages (1,2,3,4) have to be loaded from disk.

Only page 2, 3, and 4 will be modified based on redo info from log entries 11, 8 and 12.

(3):

Transaction T_1 will be rolled back. The CLR's written during undoing this update is shown below.

LSN	Type	TID	PrevLSN	UndoNxtLSN	Data
14	CLR	1	-	2	Page 2
15	CLR	1	-	1	Page 4

Part 7 Bonus: Physical Optimization (Total: 10 Bonus Points)

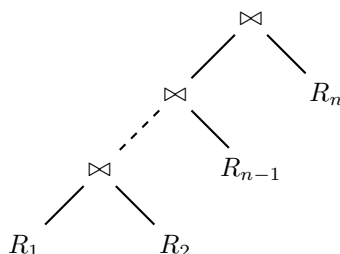
Consider the following relations $R(A, B)$, $S(C, D, E)$, $T(F, G)$ with $S = \frac{1}{10}$ (10 tuples fit on each page). The sizes and value distributions are:

$N(R) = 100$	$V(R, A) = 100$	$V(R, B) = 10$	
$N(S) = 5,000$	$V(S, C) = 500$	$V(S, D) = 2,000$	$V(S, E) = 1,000$
$N(T) = 20,000$	$V(T, F) = 1,000$	$V(T, G) = 600$	

Question 7.1 Greedy Join Enumeration (10 Points)

Use the greedy join enumeration algorithm to find the cheapest plan for the join $R \bowtie_{B=C} S \bowtie_{D=F} T$. Assume that **nested-loop** (not the block based version) is the only available join implementation with the left input being the “outer” (for each tuple from the outer we have to scan the whole inner relation). Furthermore, there are no indices defined on any of the relations (that is you have to use **sequential scan** for each of the relations). As a cost model consider the **total number of I/O operations**. For example, if you join two relations with 5,000 and 10,000 tuples with $S = \frac{1}{10}$, where the 5,000 tuple relation is the outer, then the cost would be 5,000,000 (scan the inner 5000 times) + 500 to scan the other once. The total cost is then 5,000,500 I/Os. Assume that the system supports pipelining for the outer input of a join. That is if you join the result of a join with a relation where the join result is the outer, then there is no I/O cost for scanning the outer. Also under these assumptions you never have to store join results to disk. **Hint: You will have to estimate the size of intermediate results. Use the estimation based on the number of values and not the one based on the size of the domain. Use the assumption that the number of values in a join attribute of a join result is the minimum of the number of values in the join attribute of each input.**

Write down the state after each iteration of the algorithm using the following notation. Write $((R_1, R_2), \dots, R_{n-1}), R_n)^{C, S}$ to denote a plan as shown below with I/O cost C and result size S . Alternatively you are also allowed to draw join trees as shown below.



Solution

Calculate Result Sizes:

Using the formula from class the estimated result sizes are:

$$\begin{aligned}
 T(R \bowtie S) &= \frac{T(R) \cdot T(S)}{\max(V(R, B), V(S, C))} = \frac{100 \cdot 5,000}{\max(10, 500)} = 1,000 \\
 T(S \bowtie T) &= \frac{T(S) \cdot T(T)}{\max(V(S, D), V(T, F))} = \frac{5,000 \cdot 20,000}{\max(2000, 1000)} = 50,000 \\
 T(R \bowtie T) &= T(R) \cdot T(T) = 100 \cdot 20,000 = 200,000 \\
 R(R \bowtie S \bowtie T) &= \frac{T(R) \cdot T(S) \cdot T(T)}{\max(V(R, B), V(S, C)) \cdot \max(V(S, D), V(T, F))} = \frac{100 \cdot 5,000 \cdot 20,000}{500 \cdot 2,000} = 10,000
 \end{aligned}$$

Initialization:

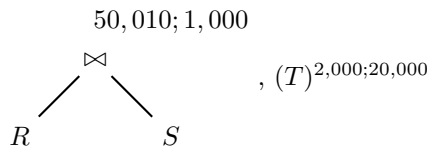
$$(R)^{10;100}, (S)^{500;5,000}, (T)^{2,000;20,000}$$

n = 1:

Here we have 6 different options how to join two of the plans from the initialization:

$$\begin{aligned}
 &(R \bowtie S)^{50,010;1,000}, (R \bowtie T)^{200,010;200,000}, (S \bowtie R)^{50,500;1,000}, \\
 &(S \bowtie T)^{10,000,500;50,000}, (T \bowtie R)^{202,000;200,000}, (T \bowtie S)^{1,002,000;50,000}
 \end{aligned}$$

As an example take the join $(R \bowtie S)$. Here R is the outer and S is the inner. The cost is computed as: For each tuple from R ($T(R)$) we have to scan S once (500 I/O). Thus, the cost is $B(R) + T(R) \cdot B(S) = 10 + 100 \cdot 500 = 50,010$ I/Os. Greedy join enumeration chooses the plan with the lowest cost $(R \bowtie S)$:



n = 2:

Now we need to consider two join options.

For $((R \bowtie S) \bowtie T)$ we pipeline the result of $(R \bowtie S)$ so the cost is:

$$Cost(R \bowtie S) + T(R \bowtie S) \cdot B(T) = 50,010 + 1,000 \cdot 2,000 = 2,050,010$$

Recall that the assumption is that only the outer input of the join can be pipelined. For $(T \bowtie (R \bowtie S))$, the result of the join $(R \bowtie S)$ is the “inner”, so we have to store the result of $R \bowtie S$ on disk resulting in $B(R \bowtie S)$ additional I/O. Since $(R \bowtie S)$ has 1,000 result tuples and $S(R \bowtie S) = S(S) + S(R) = 1/10 + 1/10 = 1/5$ it follows that $B(R \bowtie S) = 100$. Thus, the total cost is

$$Cost(R \bowtie S) + B(R \bowtie S) + B(T) + T(T) \cdot B(R \bowtie S) = 50,010 + 100 + 2,000 + 20,000 \cdot 100 = 2,052,110$$

$$(R \bowtie S) \bowtie T)^{2,050,010;10,000}, (T \bowtie (R \bowtie S))^{2,052,110;10,000}$$

