

Principles of Computer System Design

Exam

Tudor Dragan (xlq880)

January 17, 2015

Question 1: Data Processing

1. Sort-based external memory algorithm

//TODO: Insert pseudocode here

Figure 1: Sort-based external memory algorithm

For this part of the exam we must implement an algorithm that first applies an aggregated function (in our case the *count* function) on the table *friends*, and then combine the results to write up the table that consists of (uid, network, nrOfFriends). In order to achieve this, I used a modified sort-merge join algorithm.

I made the following assumptions:

1. We don't have any indices on the tables so the entries are not in sorted order in any of the two tables.
2. The friends table has only uni-directional relations, in the sense that Person1 can be friends with Person2 but it's not requested that Person2 has to be friends with Person1. (Person1 has 1 friend and Person2 has 0 friends.)
3. The table with the biggest number of records has N records.
4. We know that the main memory can hold \sqrt{N} records. If we read B records at a time in memory, the number of runs will be N/B .
5. Number of passes in Phase 2 is P then: $B(B-1)^P = N$

First I build the table with the aggregate count function for the friends table by altering the way that we merge the pages on a pass. I will explain how it is done in the following paragraphs on a concrete example:

1. Every entry in the friends table has this form : $(uid1, uid2)$. The input is split up in multiple blocks. Let us assume that we have a block that has these following $(uid1, uid2)$ relations: (3,7) (1,4) (2,3) (2,5) (1,3) (1,2) (2,4) (3,1)
2. In order to calculate the number of friends for each user, I will sort records with the form $(uid1, friendCount)$. because every record in the *friends* table essentially means 1 friend added, when reading in the first phase, the *friendCount* is 1. So we will have initially tuples with the form $(uid, 1)$ that need to be sorted by *uid*. We assume that the buffer page number is 4. So after the first phase we will have [(3,1) (1,1) (2,1) (2,1)] [(1,1) (1,1) (2,1) (3,1)].
3. The next pass will be [(1,1) (2,2) (3,1)] [(1,2) (2,1) (3,1)], Then we combine the values by adding up the number of friends by comparing the heads of the lists (the heads contain the smallest *uid*) and add up or write to disk the smallest one *uid* with the friend count.
4. Finally we will have [(1,3) (2,3) (3,2)].
5. This will continue until all the pages buffers are empty and have no more data to fill them up with.

After we create the number of friends table that has entries sorted by *uid*, we sort the *user* table. We sort the users by uid because we would like to split both sorted inputs into blocks and combine the result to (uid, network, nrOfFriends) output form. Because we have space in memory for \sqrt{N} records, each buffer block should have a size of $\sqrt{N}/2$. We compare the heads of the lists, because they always have the smallest possible id and merge the values together. We continue filling up the buffers and stop when there are no more entries to compare.

2. Hash-based external memory algorithm

//TODO: Insert pseudocode here

Figure 2: Hash-based external memory algorithm

The hash-based algorithm is very similar to the *GRACE hash join algorithm* with a slight modification that increments a counter for displaying the number of friends for each user. We make the following assumptions:

1. There are B-1 buckets, with B main memory buffers. These buckets will hold $\sqrt{N}/2$ records, where N is the highest of U and F.
2. The hashtable in the second phase is smaller than B-1 pages.

We first partition both relations U (user table) and F (friend table) via a hash function into B-1 buckets. We read one partition at a time. We iterate through it and place the record in one of the buckets by the hashing the join attribute (in our case the *uid*). Then we can be sure that if tuples of U and F join, they will wind up in corresponding buckets U_i and F_i for some *i*. After this phase we read a bucket of *U* from disk and hash it using another function and construct the hash table in memory. Then we load in a buffer for S_i and iterate through it.

When we match the *uid*, we increment the *numberOfFriends* counter in the hash table. When the buffer is empty, we write to disk the hash table with $(uid, networth, numberOfFriends)$ and load the next entries in the buffer.

3. I/O costs

Sort-merge algorithm

For the sort-merge algorithm the "expected" cost is:

$$Sort_U + Sort_F + (Pages_U + Pages_F) \quad (1)$$

Hash-based algorithm

Question 2: Distributed Transaction

1. Local wait-for-graphs

Conflict-serializability is defined by equivalence to a serial schedule (no overlapping transactions) with the same transactions, such that both schedules have the same sets of respective chronologically ordered pairs of conflicting operations (same precedence relations of respective conflicting operations).

A schedule is conflict-serializable if and only if its precedence graph has no cycles. This is a graph of nodes and vertices, where the nodes are the transaction names and the vertices are attribute collisions. The local wait-for graphs on each node are:

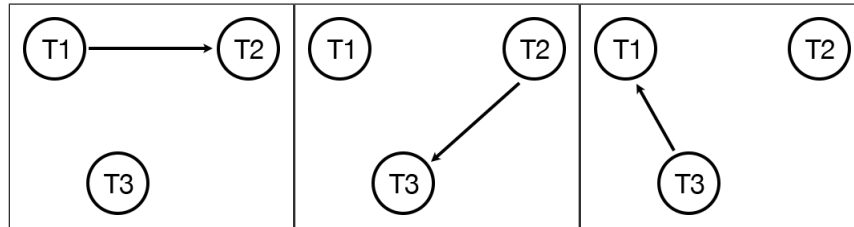


Figure 3: Local wait-for-graph for Node 1, Node 2 and Node 3

As we can see the graphs have no cycles, proving that there are no local deadlocks.

2. Global wait-for-graph

When constructing the global wait-for-graph we set up the pages that are being modified on each machine. If we have a cycle in the global graph, it means that we have conflicts and must abort the transactions. Below we have the global wait-for-graph:

When building the graph we have multiple precedences:

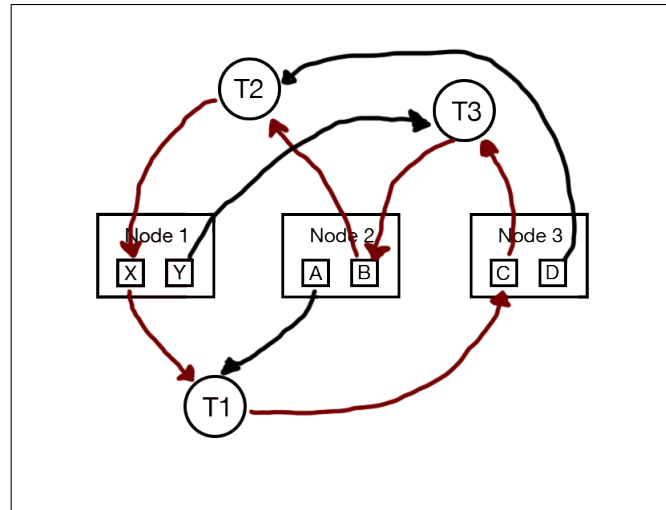


Figure 4: Global wait-for-graph

- There are $B-1$ buckets, with B main memory buffers. These buckets will hold $\sqrt{N}/2$ records, where N is the highest of U and F .
- The hashtable in the second phase is smaller than $B-1$ pages.

As we can see we have cycle between $X \rightarrow T1 \rightarrow C \rightarrow T3 \rightarrow B \rightarrow T2 \rightarrow X$. To remove the deadlock within the system we must abort modifications that are being done on X , C and B .