

CARNEGIE MELLON UNIVERSITY  
COMPUTER SCIENCE DEPARTMENT  
15-445/645 – DATABASE SYSTEMS (SPRING 2024)  
PROF. JIGNESH PATEL

Homework #5 (by Alexis Schlomer & Lan Lou) – Solutions  
Due: **Saturday Apr 20, 2024 @ 11:59pm**

**IMPORTANT:**

- Enter all of your answers into **Gradescope by 11:59pm on Saturday Apr 20, 2024.**
- **Plagiarism:** Homework may be discussed with other students, but all homework is to be completed **individually**.
- **You have to use this PDF for all of your answers.**

For your information:

- Graded out of **100** points; **4** questions total

*Revision : 2024/04/23 17:22*

Question	Points	Score
ARIES	28	
Two-Phase Commit	24	
Distributed Query Plan	18	
Miscellaneous	30	
Total:	100	

**Question 1: ARIES ..... [28 points]**

RusTub uses ARIES recovery with fuzzy checkpoints. It also has a background thread that may arbitrarily flush a dirty bufferpool page to disk at any time.

For this question, assume objects A, B, C reside in three different pages A, B, C, respectively.

LSN	WAL Record
1	<T1, BEGIN>
2	<T1, UPDATE, prev=1, B, 20→30>
3	<T2, BEGIN>
4	<T3, BEGIN>
5	<T3, UPDATE, prev=4, C, 30→40>
6	<T2, UPDATE, prev=3, A, 10→20>
7	<T2, COMMIT, prev=6>
8	<CHECKPOINT BEGIN>
9	<T1, UPDATE, prev=2, A, 20→30>
10	<CHECKPOINT END, ATT={T1, T2, T3}, DPT={C}>
11	<T3, UPDATE, prev=5, C, 40→50>
12	<T2, TXN-END>
13	<CHECKPOINT BEGIN>
14	<T1, COMMIT, prev=9>
15	<T4, BEGIN>
16	<CHECKPOINT END, ATT={T1, T3}, DPT={?}>
17	<T4, UPDATE, prev=15, B, 30→20>
18	<T4, UPDATE, prev=17, A, 30→20>
19	<T4, ABORT, prev=18>
20	<T4, CLR, prev=19, A, 20→30, undoNext=17>

Figure 1: WAL

- (a) Suppose the system crashes and, when it recovers, the WAL contains the first 10 records (up to <CHECKPOINT END, ATT={T1, T2, T3}, DPT={C}>). Of the object states below, which states are possibly stored on disk before recovery starts? Select all that apply.

- i. [4 points]  A=10     A=20     A=30  
 ii. [4 points]  B=20     B=30  
 iii. [4 points]  C=30     C=40

**Solution:** Page A is not logged as dirty in the checkpoint record, so the value can't be 10. But LSN 9 changes page A, so it can be either 20 or 30. Page B is not logged as dirty in the checkpoint record, and there is no followup updates between LSN 8 and LSN 10, so its value should be 30. Page C is logged as dirty in the checkpoint. It's uncertain whether it hasn't been flushed at all or if it has been flushed at some stage, so its value can be 30 or 40. Pages on disk will not contain any updates later than the last flushed log record.

(b) [4 points] Select all possible values of DPT in record 16.

- A    B    C    A, B    A, C    B, C    A, B, C    None of them

**Solution:** Dirty pages can be flushed to disks at any time arbitrarily. But LSN 10 ensures that the value stored in disk for Page B has already been updated from 20 to 30, and there is no followup updates, so B cannot appear in the DPT. Both Page A and C may be present in the DPT due to the following updates after LSN 8 (not to mention C has already been logged dirty at the last checkpoint).

(c) [4 points] For next 3 questions, assume that the database restarts and finds all log records up to LSN 20 in the WAL. Also assume the DPT is {C} for LSN 16. According to the lecture, which pages the analysis phase may select to be redone? Select all that apply.

- A    B    C    None of them

**Solution:** We redo all pages in the DPT and those modified after LSN 13.

(d) [4 points] Select all transactions that should be undone during recovery.

- T1    T2    T3    T4    None of them

**Solution:** All uncommitted or explicitly aborted transactions should be undone.

(e) [4 points] How many new CLR records will be appended to the WAL after the database fully recovers?

- 0    1    2    3    4    5    6

**Solution:** T3 will produce 2 more CLR records. T4 will produce 1 more CLR record.

**Question 2: Two-Phase Commit..... [24 points]**

Consider a distributed transaction  $T$  operating under the two-phase commit protocol with the *early acknowledgement optimization*. Let  $N_0$  be the *coordinator* node, and  $N_1, N_2, N_3$  be the *participant* nodes. Let the  $C$  be the client application issuing the transaction request. Assume that the client communicates directly with the coordinator node.

The following messages have been sent:

time	message
1	$C$ to $N_0$ : "REQUEST: COMMIT"
2	$N_0$ to $N_2$ : "Phase1: PREPARE"
3	$N_0$ to $N_3$ : "Phase1: PREPARE"
4	$N_2$ to $N_0$ : "OK"
5	$N_0$ to $N_1$ : "Phase1: PREPARE"
6	$N_1$ to $N_0$ : "OK"
7	$N_3$ to $N_0$ : "OK"

Figure 2: Two-Phase Commit messages for transaction  $T$

(a) [6 points] Who should send message(s) next at time 8 in Figure 2? Select *all* the possible answers.

- $C$
- $N_0$
- $N_1$
- $N_2$
- $N_3$
- It is not possible to determine

**Solution:** At the time when all participant nodes have responded with "OK",  $T$  is considered to be committed.  $N_0$  must notify the participant nodes  $N_1, N_2$ , and  $N_3$  of the decision to commit. At this time  $N_0$  should also notify the client application since the protocol is running with the early acknowledgement optimization.

(b) [6 points] To whom? Again, select *all* the possible answers.

- $C$
- $N_0$
- $N_1$
- $N_2$
- $N_3$
- It is not possible to determine

**Solution:** At the time when all participant nodes have responded with "OK",  $T$  is considered to be committed.  $N_0$  must notify the participant nodes  $N_1, N_2$ , and  $N_3$  of the decision to commit. At this time  $N_0$  should also notify the client application since the protocol is running with the early acknowledgement optimization.

- (c) [6 points] Suppose that  $N_0$  received the “**ABORT**” response from  $N_3$  at time 7 in Figure 2. What should happen under the two-phase commit protocol in this scenario?
- $N_0$  resends “**Phase1 : PREPARE**” to  $N_2$
  - $N_1$  resends “**OK**” to  $N_0$
  - $N_0$  sends “**Phase2 : COMMIT**” all of the participant nodes
  - $N_0$  sends “**ABORT**” all of the participant nodes
  - $N_0$  resends “**Phase1 : PREPARE**” to all of the participant nodes
  - It is not possible to determine

**Solution:** The coordinator ( $N_0$ ) will mark the transaction as aborted. 2PC requires that *all* participants respond with “**OK**”.

- (d) [6 points] Suppose that  $N_0$  successfully receives all of the “**OK**” messages from the participants from the first phase. It then sends the “**Phase2 : COMMIT**” message to all of the participants but  $N_1$  and  $N_3$  crash before they receive this message. What is the status of the transaction  $T$  when  $N_1$  comes back on-line?
- $T$ 's status is *aborted*
  - $T$ 's status is *committed*
  - It is not possible to determine

**Solution:** Once the coordinator ( $N_0$ ) gets a “**OK**” message from *all* participants, then the transaction is considered to be committed even though a node may crash during the second phase. In this example,  $N_1$  and  $N_3$  would have to restore  $T$  when it comes back on-line.

**Question 3: Distributed Query Plan ..... [18 points]**

The CMU-DB team is optimizing distributed databases for aggregations and joins in their new project, AutoLoo. This project aims to enable efficient computation of distributed joins followed by aggregations.

Given the following schema:

```
CREATE TABLE cust (PRIMARY KEY cust_id int, name VARCHAR, loc_id int);
CREATE TABLE txn (PRIMARY KEY txn_id int, cust_id int, amount DECIMAL, year
CREATE TABLE loc (PRIMARY KEY loc_id int, region_name VARCHAR);
```

AutoLoo partitions these tables across nodes based on the partition key. Consider the query for analyzing total spending per customer region for the year 2024:

```
SELECT l.region_name, SUM(t.amount) AS total_spending
FROM txn t
JOIN cust c ON t.cust_id = c.cust_id
JOIN loc l ON c.loc_id = l.loc_id
WHERE t.year = 2024
GROUP BY l.region_name;
```

You can make the following assumptions:

1. There are 5 nodes in the system.
  2. The customer table contains 20,000 rows.
  3. The transaction table contains 50,000 rows for the year 2024.
  4. The location table contains 500 rows.
- (a) [5 points] Which data distribution strategy minimizes the total network data transfer for the given query?
- A) Partition customer and transaction tables by customer\_id range, and replicate location table across all nodes.
  - B) Partition transaction table by transaction\_id range, and customer and location tables by location\_id range.
  - C) Replicate customer and location tables across all nodes, and partition transaction table by customer\_id.
  - D) Partition all tables randomly without any specific range or replication strategy.

**Solution:** This strategy minimizes cross-node data transfer by localizing joins involving customer and transaction tables to individual nodes and avoids the need to transfer location data across nodes because it is replicated everywhere.

- (b) **[5 points]** Assuming the selected strategy from question (a) is implemented, what is the estimated total data transferred over the network for the **join** operation? **Hint:** Only the central node can perform aggregations.
- A) Less than 5,000 rows
  - B) Between 5,001 to 25,000 rows
  - C) **Between 25,001 to 100,000 rows**
  - D) More than 100,000 rows

**Solution:** Given that location data is replicated across all nodes and joins are localized, the data transfer will be the result of all local joins. Worst case scenario  $10,000 \times 5 = 50,000$  rows get transmitted to the central node, which is simply the max cardinality of the join output.

- (c) **[5 points]** If AutoLoo introduces a feature that allows for intermediate aggregation results to be computed on each node before being sent to a central node for final aggregation, how will this impact the total network data transfer?
- A) **Less than 5,000 rows**
  - B) Between 5,001 to 25,000 rows
  - C) Between 25,001 to 100,000 rows
  - D) More than 100,000 rows

**Solution:** Here, the only data transfer will be the result of local aggregations needing to be combined at a central node for final computation. Since we have 5 nodes and a maximum of 500 different locations, this translates into an upper bound of  $5 \times 500 = 2,500$  row transfers.

- (d) **[3 points]** What are the primary drawbacks of implementing a feature that allows for intermediate aggregation results to be computed on each node before sending these results to a central node for final aggregation? Consider the impact on system resources. Select all that apply.
- A) It significantly increases the amount of data transferred over the network.
  - B) **It increases the computational load on each node.**
  - C) **It increases the memory usage on each node due to the storage of intermediate results.**
  - D) It decreases the overall system performance.

**Solution:** Performing intermediate aggregations locally increases each node's computational load and memory usage, as they must handle and store aggregation results. This method reduces network traffic but demands more computational and memory resources. Option A is incorrect since the goal is to decrease data transfer. Option D is too general; local aggregations often enhance performance by easing network demands.

**Question 4: Miscellaneous ..... [30 points]**

- (a) [3 points] A distributed DBMS can commit transactions and automatically compensate for network partitioning without any loss of data consistency.

- True  
 False

**Solution:** This question was ill-phrased and meant to be: "A distributed DBMS can immediately commit a transaction under a network partitioning without any loss of data consistency." The intended answer to that is False, as this is a violation of the CAP theorem.

- (b) [3 points] ARIES employs a single pass of the log during the recovery process to handle both redo and undo operations.

- True  
 False

**Solution:** First there is a REDO phase, then an UNDO phase.

- (c) [3 points] The CAP theorem implies that a distributed system cannot simultaneously guarantee consistency, availability, and partition tolerance.

- True  
 False

**Solution:** By definition.

- (d) [3 points] In ARIES, only transactions that commit will have an associated "TXN-END" record in the log.

- True  
 False

**Solution:** Transactions will commit when the COMMIT log is written, TXN-END is an optimization in the ARIES protocol that indicates it can be removed from the transaction table.

- (e) [3 points] In the context of distributed DBMS, data replication increases availability but can lead to challenges in maintaining data consistency across nodes.

- True  
 False



**Solution:** CAP limitations.

- (f) [3 points] Both PAXOS and Two-Phase Commit protocols can be used to implement distributed transactions.

True

False

**Solution:** 2PC is a degenerate case of Paxos, which is a consensus protocol for distributed transactions.

- (g) [3 points] In reference to recovery algorithms that use a write-ahead log (WAL). Under NO-STEAL + FORCE policy, a DBMS will have to undo the changes of an aborted transaction during recovery.

True

False

**Solution:** The FORCE policy guarantees that all changes from committed transactions are immediately written to disk, ensuring durability, while it is the NO-STEAL policy that prevents changes by uncommitted transactions from being written to disk, avoiding the need for undo operations on disk for those transactions.

- (h) [3 points] Fuzzy checkpoints need to block the execution of all transactions while a consistent snapshot is written to disk.

True

False

**Solution:** See motivation of checkpointing.

- (i) [3 points] With consistent hashing, if a node fails then all keys must be reshuffled among the remaining nodes.

True

False

**Solution:** See motivation of consistent hashing.

- (j) [3 points] In a system with strong consistency requirements, it is best for the DBMS to implement active-passive replication with asynchronous replication and continuous log streaming.

True

False

**Solution:** Asynchronous replication cannot guarantee consistency (see CAP).