

Project 4: Logging and Recovery

Due April 17 at 11:55:00 pm

Introduction

It's a sad fact of life: computers crash. One minute everything is working just fine, and the next minute everything that was in volatile memory is gone forever. To deal with this problem on your personal machine, you've probably gotten into the habit of saving frequently. In this project, you will implement an algorithm that database systems commonly use to address the same problem without slowing down the system too much.

Some Basics of Transactions¹

A *transaction* is any one execution of a user program in a database management system (DBMS). For example, suppose Ada and Bob both have accounts at the same bank. Ada logs onto the bank's website and instructs the bank to make a payment of \$100 from her account into Bob's account. If the bank's computer runs normally, it will deduct \$100 from Ada's account balance, add \$100 to Bob's account, and trigger emails to Ada and Bob confirming the transfer. Those steps—deducting, adding, and confirming that it was done—make up one transaction.

But suppose that the moment after the bank's computer deducts the money from Ada's account balance, it crashes. When it reboots, the instructions about that \$100 that were in its memory are gone. We would like to ensure that if this happens, the \$100 has not simply vanished—it is back in Ada's account, or it is in Bob's account. Either the whole transaction goes through, or the part of the transaction that happened before the crash gets undone. Another way of saying this is that we want a transaction to be **atomic**.

Suppose the timing of the crash was a little different. The bank's computer read Ada's balance into memory and modified it there, read Bob's balance into memory and modified it there, and triggered the emails. Then, just as it was about to write the updated balances back onto disk, it crashed. We want to make sure that if the system says a transaction completed successfully, the changes the transaction made will persist even through a crash. We call this **durability**.

(In lecture, you will also learn about two other properties that transactions should have, called **consistency** and **isolation**. These have to do with concurrency and are important to understand, but you don't need to worry about them for this project.)

To ensure atomicity and durability while minimizing the harm to performance, DBMSs use **write-ahead logging** and with a recovery algorithm, such as **ARIES**. In this project, you will be implementing a single-threaded version of ARIES.

¹ Transactions will be covered in lecture soon, but if you just can't wait to learn more, check out chapter 16 of the textbook.

What you need to do:

Read.

Read sections 16.7 through 16.7.2 and 18.1 through 18.6 in your textbook. **This is important.** The rest of this project will not make sense until you have understood this material.

(Since we have been giving you reading assignments all semester, we assume you already own the textbook or know that it can be checked out for two hours at a time from the Art Architecture & Engineering library's course reserve desk on the second floor.)

Download.

Download the recovery simulator from CTools and unzip it. You will see that the eecs484p4 directory contains directories called StorageEngine, StudentComponent, output, correct, and testcases. The section on *Understanding the Simulator* below explains what these pieces do. You'll also want to spend some time familiarizing yourself with the header files.

Implement.

Implement ARIES in LogMgr.cpp.

Do not alter the other files in any way. The autograder will compile and run your LogMgr.cpp with unaltered copies of all of the other files, so it is in your best interest to test with exactly the same files. Please do not add any header files; if you need helper functions, you may declare them and define them in LogMgr.cpp.

You may use the C++11 STL, **except** that you may **not** use **anything that reads from or writes to disk**.

Submit.

An autograder will be available by Saturday, April 4. The URL will be announced when the autograder is open. Submit your LogMgr.cpp file to the autograder. You will need your UMich credentials to log in. Enter your partner's unickname in the box, or leave it blank if you are working alone.

Each unickname may submit up to three times per day. If you work with a partner, this means the pair of you can have a total of six submissions per day. Please take advantage of this, and start submitting early and often.

We will grade the last submission you or your partner turn in before the end of the late days (11:55 pm on April 21). If you submit before the regular deadline and also during the late days, we will grade your late submission, and you will incur the late penalty.

Don't forget to submit a partner note as usual through the CTools assignment!

Grading

Your grade will be based on your code's performance on the autograder test cases. We will diff your log file and your recovered db file against our known correct files for each test case. For each test case, you earn 5 points for generating the correct log file and 5 points for generating the correct db file.

The test cases on the autograder are different from the public test cases we have provided you. Passing the public test cases means that you are on the right path, but doesn't ensure full marks on the private test case. The autograder test cases may also cover corner cases the public test cases do not. Thus, make sure you think carefully about all the cases that the algorithm needs to handle.

Understanding the Simulator

The simulator has two main pieces: the StorageEngine and the LogMgr. The StorageEngine handles all disk reads and writes. The LogMgr handles logging and recovery. If the LogMgr needs to write something to disk, it must go through the StorageEngine. **Your implementation of LogMgr.cpp must not include any static variables or functions, and it must not write directly to disk. If LogMgr wants to write to disk, it must use StorageEngine's public methods.**

The StorageEngine keeps track of Pages. Each Page has a page_id, a pageLSN, a dirty bit, and a string of data.² All Pages are kept on disk. If a page needs to be read or written, StorageEngine adds it to the in-memory page buffer. If the page buffer is full, the StorageEngine chooses a page to flush to disk, informs the LogMgr of the flush through LogMgr's pageFlushed function (Hint: that probably means LogMgr should do something when the page is flushed), and flushes it, freeing up memory for the next page.

Sometimes LogMgr will need to write Pages; for example, recovering from a crash and aborting a transaction both require LogMgr to alter Pages. When this is necessary, LogMgr can call StorageEngine::pageWrite. You should always try to minimize the number of Pages LogMgr writes. To force you to do this, StorageEngine limits the number of pageWrites you are permitted.³ If you try to exceed the limit, pageWrites will stop writing pages and return false. If this happens, LogMgr should stop what it is doing.

The LogMgr also needs to keep the log. A log is made up of LogRecords of various types. You can keep the most recent part of the log (called the log tail) in memory, but sometimes you will need to write the log to disk. LogMgr can call a LogRecord's toString method to transform the

² In a real database, of course, pages contain records, but here we've abstracted them away, so you don't need to worry about them.

³ The limit at any given time is set by the testcase. Sometimes we give you enough page writes to finish an operation. Sometimes, to simulate a crash in the middle of an operation, you will run out of pageWrites early, and the next event in the testcase will be a crash.

LogRecord into a string, and then pass a string to StorageEngine::updateLog to append a string to the log on disk. The log on disk will have one record per line; you can append multi-line strings to it if you want to add more than one record at once. If LogMgr wants to read old log entries, it can call StorageEngine::getLog(), which will return a (multi-line) string. LogRecord::stringToRecordPtr can parse a line of that string and give you a pointer to a LogRecord of that line.

StorageEngine provides a few other utilities LogMgr might find useful:

- StorageEngine::nextLSN provides a unique id number assigned in monotonically increasing order
- StorageEngine::store_master and StorageEngine::get_master write an int to stable storage and retrieve it.
- StorageEngine::getLSN(int page_id) returns the page LSN of the specified page.

In the StorageEngine directory, you will also find main.cpp. Main takes the relative path to a testcase as a command line argument. For example, once you have compiled the whole project, you might call `./main.o testcases/test00` from `eeecs484p4` to run testcase 0. Main calls the `runTestcase` function, which will read in the testcase line by line and run a simulation.

Test Cases

The first line of the testcase specifies a "database" file (db) to use for this run. For this simulation, we are using an oversimplified format, so that you can see the effects of writing to pages: a database file is a text file, and each line of text represents a page of data on disk. If the first line of a database file looks like this:

```
-1 xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
```

it means that the page with `page_id` 1 (because it is the first line of the file) has `pageLSN` -1, and the data currently written to that page is 50 consecutive x's.

The remaining lines of the testcase are instructions to write, abort, commit, make a checkpoint, crash, or end the simulation. Each of these is described in more detail below, with examples.

- `write` (e.g., `1 write 34 27 ABC`): An instruction from a particular transaction (1) to write an update (ABC) to a particular page (34) at a particular offset (27). Main's `runTestcase` function calls StorageEngine's `write` function, which in turn calls LogMgr's `write` function, giving it the transaction id, page id, offset, input string, and what was already written at that page-offset location. StorageEngine expects to get back a `pageLSN` for the write; once it receives this, it updates the page.
- `abort` (e.g., `1 abort 5`): An instruction to abort a particular transaction (1), and a cap on the number of page writes the LogMgr is permitted (5). `RunTestcase` calls StorageEngine's `abort` function; StorageEngine sets the number of page writes permitted to five and then calls `LogMgr::abort(1)`.
- `commit` (e.g., `1 commit`): An instruction to commit a transaction (1). `RunTestcase` calls `LogMgr::commit(1)`.

- `checkpoint (checkpoint)`: An instruction to create a checkpoint. `RunTestcase` calls `LogMgr::checkpoint()`.
- `crash (e.g., crash {2 5})`: An instruction to simulate one or more crashes and recovery, with the specified number of writes permitted during each recovery phase. In this example, the first crash happens, and the `LogMgr` will be allowed two writes before another crash happens; `LogMgr` gets five writes to recover. Sometimes the number of writes will allow a full recovery before the next crash; sometimes it will not, and the second crash will begin with `StorageEngine` no longer writing when `LogMgr` requests it. When crashing, `RunTestcase` will destroy the current instance of `LogMgr` and create a new one, then call `StorageEngine::crash()`. `StorageEngine` will destroy all of the `Pages` it has in memory, then call `LogMgr`'s `recover` method. `LogMgr` is responsible for recovering according to the ARIES protocol.
- `end (end)`: An instruction to end the simulation. Saves the final state of the on-disk pages and exits the program.

At the end of the testcase, the simulator will have generated two files: a log and a modified version of the input db file. These can be found in the `output/log/` and `output/dbs/`, respectively. Each filename includes the number of the corresponding testcase.

NOTE: If you already have a file in the log folder with the same name as the new one, the new log entries will be appended to the file, rather than overwriting it. You might want to rename or delete old log files to avoid some confusing bugs.

Reading the Log

The log file is saved as a multi-line string, with one line for each log entry, so that you can easily review it. The format is loosely based on Figure 18.2 in the textbook. For example, a log might have content like this (heading and table borders added for clarity):

LSN	prevLSN	tx_id	type	page_id	offset	before image	after image
2	-1	1	update	5	0	xxxxxxx	update10
3	-1	2	update	3	0	xxxxxxx	update20
4	3	2	commit				

End checkpoint log records also include a record of the transaction table and the dirty page table. Each table is enclosed in curly braces, and each entry in the table is enclosed in square braces. For example,

LSN	prevLSN	tx_id	type	TX_table	Dirty Page Table
5	4	-1	end_checkpoint	{ [1 2 U] [2 3 U] }	{ [3 3] [5 2] }

This transaction table has two entries, and so does the dirty page table.

We have given you five public test cases. The correct output for these can be found in the `correct/logs/` and `correct/dbs/` directories.

Trouble Debugging?

If you are using gdb on the CAEN computers to debug, you might find that it has some difficulty displaying `std::map` and other data structures that your implementation might use. To fix this, you can use Pretty Print:

Make a folder for it wherever you like to keep such things (mine is `~/gdb_printers`) and, from inside that folder, run

```
svn co svn://gcc.gnu.org/svn/gcc/branches/gcc-4_6-branch/libstdc++-v3/python
```

Once you have that, create a file `~/gdbinit` (or append to it if the file already exists) and add the following lines to it:

```
python
import sys
sys.path.insert(0, '/path/to/gdb_printers/python')
from libstdcxx.v6.printers import register_libstdcxx_printers
register_libstdcxx_printers (None)
end
```

Make sure you replace `/path/to/gdb_printers/python` with the actual path.

Now try running `gdb --args ./main.o testcases/test00` again, and you should be able to view the maps properly.

Important Information about the Honor Code

Your assignment is to implement ARIES within the simulator environment. Because you are smart and good at hacking, you might be tempted to use some other method to generate correct output files. Don't try it.

Trying to trick the autograder into giving you a good grade when you have not actually done the assigned work is cheating. We have technical measures built into the autograder to detect all of the obvious hacks (because we are smart and good at hacking, too), and we will review code by hand. If you are caught, you will receive a zero on the project and will be referred to the honor council.

As mentioned above, **your `LogMgr.cpp` file must not attempt to read or write any file directly, and it must not include any static functions or variables. Any attempt to do so will be detected and logged as signs of cheating.** While these are not the only signs we are looking for, we wanted to warn you about these in particular, since many of you might be used

to using this functionality with no ill intent. If you have any questions about whether a strategy you plan to use would be construed as cheating, please ask the course staff before using it.

A Few Clarifications:

- When you're undoing a transaction, any time you see a log record for that transaction, you need to get the prevLSN from that log record and add it to your toUndo list (or end the transaction if the prevLSN is -1). If the record is a CLR or an update record, you then process it the way the book says; if it's something else (like an abort record), you just move on to the next item in your toUndo list.
- You may have noticed in lecture and some of the problems in the textbook that the log tail is sometimes flushed to disk without a triggering commit or page write or end checkpoint. ARIES permits a regular background flushing of the log tail; for example, a certain amount of space in memory called a log buffer may be set aside, and any time it fills up, the log tail can be flushed to disk. **You don't need to worry about that background flushing for this project.** Any reasonably sized log buffer would be bigger than the log tail ever gets in these test cases.
- The book's description of abort is a little brief. What you need to know is that your abort function should write an abort record, then do basically the same thing as undo, only on a smaller scale. For undo, you need to look at all of the loser transactions and put their most recent LSNs in your toUndo queue. For abort, you just need to find the most recent LSN for one transaction you are aborting. Look at the log record for that LSN. Its prevLSN field will tell you which log record you need to look at next (or if you're done); its type will tell you if you need to process it the way you process records during undo. When you get to a log record with prevLSN = -1, process that record, then write your end log record, and you're done.