15-440/15-640 Distributed Systems

Fault Tolerance, Logging and Recovery

Lecture #11, Tuesday Oct 5th 2021

Announcements

- HW2 Released. Due 10/10, <u>No Late Days</u>
- P1: Part A Due (10/8), Part B (10/14)
 - Note, no Project OH on the due dates of Part A and Part B
- Regrade Requests: Projects, Homeworks and Midterms
 - Within 1 week of the scores/grades being posted on Gradescope
 - Except for when mitigated by university grade deadlines
- Midterm-1 In Class, During Class Time
 - Tuesday 10/12, 10:10 11:30
 - Please try and arrive early, 10:00-10:05am, get settled.
 - Topics covered, everything in the 1st half of class

Today's Lecture Outline

- Motivation Fault Tolerance
- Fault Tolerance using Checkpoints
- Fault Tolerance using Logging and Recovery
- Logging and Recovery in Practice: ARIES

What is Fault Tolerance?

Dealing successfully with partial failure within a distributed system

- System that cannot meet its goals => faults
- Recover from all kinds of faults:
 - **Transient**: appears once, then disappears
 - Intermittent: occurs, vanishes, reappears
 - **Permanent**: requires replacement / repair

Dependability

Dependability implies the following:

- 1. Availability
- 2. Reliability
- 3. Safety
- 4. Maintainability

Dependability Concepts

- *Availability* the system is ready to be used immediately.
 - "High availability": system is ready at any given time, with high probability.
- *Reliability* the system runs continuously without failure.
 - "High reliability": system works without interruption during a long period of time.

Subtle difference. Consider:

- Random but rare failures (one millisecond every hour)
- Predictable maintenance (two weeks every year)

Dependability Concepts

- *Safety* if a system fails, nothing catastrophic will happen
- Maintainability when a system fails, it can be repaired easily and quickly (sometimes, without its users noticing the failure). Also called Recovery.

Failure Models

Type of Failure	Description
Crash failure	Server halts, working correctly before.
Omission failure Receive omission Send omission	Server fails to respond to request Server fails to receive incoming msg Server fails to send msg
Timing failure	Server's response outside specified time interval
Response failure Value failure State transition failure	Incorrect server response Wrong value of response Deviation from flow of control
Arbitrary (Byzantine) failure	Server may produce arbitrary responses at arbitrary times

Masking Failures by Redundancy

- Information Redundancy add extra bits to allow for error detection/recovery (Hamming codes: detect 2-bit errors, correct 1-bit errors)
- *Time Redundancy* perform operation and, if need be, perform it again.
 (Purpose of transactions: BEGIN/END/COMMIT/ABORT)
- 3. Physical Redundancy add extra (duplicate) hardware and/or software to the system.

Can you think of Physical redundancy in Nature?

Redundancy in Electronics



Redundancy is Expensive







 But without redundancy, we need to recover after a crash.

Recovery Strategies

When a failure occurs, we need to bring the system into an error free state (recovery).

- **1. Backward Recovery**: return the system to some previous correct state (using *checkpoints*), then continue executing.
- 2. Forward Recovery: bring the system into a correct new state, from which it can then continue to execute.
 - Erasure coding \rightarrow Forward Error Correction

Forward and Backward Recovery

- Major disadvantage of Backward Recovery:
 - Checkpointing can be very expensive (especially when errors are very rare).
- Major disadvantage of Forward Recovery:
 - In order to work, all potential errors need to be accounted for *up-front*.
 - "Harder": the recovery mechanism needs to know how do to bring the system *forward* to a correct state.
 - Clever algorithms and mathematical structures help

Today's Lecture Outline

- Motivation Fault Tolerance
- Fault Tolerance using Checkpoints
- Fault Tolerance using Logging and Recovery
- Logging and Recovery in Practice: ARIES

Checkpointing

- Checkpoint: snapshot the state of the DS
 - Complete state including messages received / sent
- Checkpointing for fault tolerance:
 - Periodically save system state
 - Stored in reliable storage that can withstand targeted failure
 - Roll back to error-free state in case of failure and re-execute
- Rollback upon failure
 - Restore state to that of last checkpoint
 - All intervening computation wasted

Checkpointing assumes– Stable Storage



(a) Stable storage.(b) Crash after drive 1 is updated.(c) Bad spot.

Checkpointing

Design decisions and tradeoffs:

- How often to checkpoint?
 - Frequent checkpoints: pros and cons?
 - Pros: Less roll back (and hence wasted work/time)
 - Cons: Overhead during normal operation (Significant IO traffic, time overhead,..)
- Independent or synchronous?

Independent Checkpointing



A recovery line to detect the correct distributed snapshot This becomes challenging if checkpoints are un-coordinated

The Domino Effect



The domino effect – Cascaded rollback

P2 crashes, roll back, but 2 checkpoints inconsistent (P2 shows m received, but P1 does not show m sent)

Coordinated Checkpointing

- Key idea: each process takes a checkpoint after a globally coordinated action. (Why?)
- Simple Solution: 2-phase blocking protocol
 - Coordinator multicast *checkpoint_REQUEST* message
 - Participants receive message, takes a checkpoint, stops sending (application) messages and queues them, and sends back checkpoint_ACK
 - Once all participants ACK, coordinator sends *checkpoint_DONE* to allow blocked processes to go on
- Optimization: consider only processes that depend on the recovery of the coordinator (those it sent a message since last checkpoint)

Successful Coordinated Checkpoint



Unsuccessful Coordinated Checkpoint



Checkpoints can fail, if participant sent msg before *checkpoint_REQUEST* and receiver gets msg after *checkpoint_REQUEST*. Then: abort and do try another coordinated checkpoint soon.

Today's Lecture Outline

- Motivation Fault Tolerance
- Fault Tolerance using Checkpoints
- Fault Tolerance using Logging and Recovery
- Logging and Recovery in Practice: ARIES

Goal: Make transactions Reliable

- ... in the presence of failures
 - Machines can crash: disk contents (OK), memory (volatile)
 - Assume that machines don't misbehave
 - Networks are flaky, packet loss, handle using timeouts
- If we store database state in memory, a crash will cause loss of "Durability".
- May violate atomicity, i.e. recover such that uncommited transactions COMMIT or ABORT.
- General idea: store enough information to disk to determine global state (in the form of a LOG)

Challenges:

- Disk performance is poor (vs memory)
 - Cannot save all transactions to disk
 - Memory typically several orders of magnitude faster
- Same general idea: store enough data on disk so as to recover to a valid state after a crash:
 - Shadow pages and Write-ahead Logging (WAL)
 - Idea is to provide Atomicity and Durability

Shadow Paging

- Provide Atomicity and Durability
- "page" = unit of storage
- Idea: When writing a page, make a "shadow" copy
 - No references from other pages, edit easily!
- ABORT: discard shadow page
- COMMIT: Make shadow page "real". Update pointers to data on this page from other pages
 - Can be done atomically
- Essentially "copy-on-write" to avoid in-place page update

Write-Ahead Logging (WAL)

- Provide Atomicity and Durability
- Idea:
 - Create a log recording every update to database
 - Updated pages are kept in memory (page cache)
 - Log written to disk before the corresponding page is written to disk
 - Logs typically store both REDO and UNDO operations
 - Updates considered reliable when log entry stored on disk
 - After a crash, recover by replaying log entries to reconstruct correct state
- WAL is more common
 - Pro: Fewer random disk operations

Today's Lecture Outline

- Motivation Fault Tolerance
- Fault Tolerance using Checkpoints
- Fault Tolerance using Logging and Recovery
- Logging and Recovery in Practice: ARIES

ARIES Recovery Algorithms

- ARIES: Algorithms for Recovery and Isolation Exploiting Semantics
- Used in major databases
 - IBM DB2 and Microsoft SQL Server
 - Deals with many practical issues
- Principles
 - Write-ahead logging
 - Repeating history during Redo
 - Logging changes during Undo

Simple mental model for the database:

- fixed size PAGES, storage at page level
- Pages on disk, some also in memory (page cache)
- "Dirty pages": page in memory differs from one on disk
- Transaction Table (TT): All active TXNS
- Dirty Page Table (DPT): all dirty pages in memory

- Log: View as sequence of entries, sequential number
 - Log-Sequence Number (LSN)
- Goal: Reconstruct global consistent state using
 - Log files + disk contents + (page cache)
- Logs consist of sequence of records
 - What do we need to log?
 - Seq#, which transaction, operation type, what changed...

- Logs consist of sequence of records
 - LSN: [prevLSN, TID, type, pageID, newVal, oldVal]
 - PrevLSN forms a backward chain of operations for each TID
 - Allows REDO operations to bring a page up to date and UNDO an update reverting to an earlier version (why?)
- Transaction Table (TT): All TXNS not written to disk
 - Includes lastLSN: Seq Num of the last log entry they caused
- Dirty Page Table (DPT): all dirty pages in memory
 - Includes recoveryLSN: first log entry to make page dirty
- Each page has a pageLSN: latest log entry applied to that page
 - When in memory and on disk
 - (Denoted as just LSN in the following example)

ARIES (WAL): Data Structures

TT: Transaction Table

TID	LastLSN
1	567
2	42
7	67
8	12

TID: Transaction ID

LastLSN: LSN of the most recent log record seen for this Transaction. i.e. latest change

DPT: Dirty Page Table

pageID	recoveryLSN
42	567
46	568
77	34
3	42

pageID: key/ID of a page

recoveryLSN: LSN of first log record that made page dirty i.e. earliest change to page









TT	
TID	LastLSN
1	4
2	3

DPT	
pageID	recoveryLSN
42	1
46	3

- Commit a transaction
 - Log records written to disk until commit entry
 - No need to update actual disk pages, log file has information
 - Keep "tail" of log file in memory => not commits
 - If the tail gets wiped out (crash), then partially executed transactions will be lost. Can still recover to reliable state
- Abort a transaction
 - Locate last entry from TT, undo all updates so far
 - Use PrevLSN to revert updates to pages to start of TXN
 - A special log entry added for each undo as well (why?)
 - Type = Compensation Log Record (CLR)

Recovery using WAL – 3 passes

- 1. Analysis Pass
 - Begin from Start or last checkpoint
 - Identify the dirty pages and active txns at the time of the crash
 - Maintain a list of txns and remove them from TT if committed
 - At the end of the analysis phase has
 - TT: all active transactions at the time of the crash
 - DPT: dirty pages along with recoveryLSN

Recovery using WAL – 3 passes

- 2. Recovery Pass (redo pass)
 - Replay log forward (i.e., repeat actions) from an appropriate point
 - Bring everything to a state at the time of the crash
- 3. Undo Pass
 - Replay log file backward, revert any changes made by transactions that had not committed (use PrevLSN)
 - For each write Compensation Log Record (CLR)
 - Once reach entry without PrevLSN \rightarrow done
 - At the end: all updates from committed txns reapplied and all updates from uncommitted txns reverted

ARIES (WAL): Same example

TT: Transaction Table

TID	LastLSN
1	567
2	42
7	67
8	12

TID: Transaction ID

LastLSN: LSN of the most recent log record seen for this Transaction. i.e. latest change

DPT: Dirty Page Table

pageID	recoveryLSN
42	567
46	568
77	34
3	42

pageID: key/ID of a page

recoveryLSN: LSN of first log record that made page dirty i.e. earliest change to page









TT	
TID	LastLSN
1	4
2	3

DPT	
pageID	recoveryLSN
42	1
46	3





2PC works great with WAL/ARIES

- WAL can integrate with 2PC
 - Have additional log entries that capture 2PC operation
 - **Coordinator:** Include list of participants
 - Participant: Indicates coordinator
 - Votes to commit or abort
 - Indication from coordinator to Commit/Abort

Optimizing WAL

- As described earlier:
 - Replay operations back to the beginning of time
 - Log file would be kept forever (\rightarrow entire Database)
- In practice, we can do better with CHECKPOINT
 - Periodically save DPT, TT
 - Store any dirty pages to disk, indicate in LOG file
 - Prune initial portion of log file: All transactions upto checkpoint have been committed or aborted.

Summary

- Basic concepts for Fault Tolerant Systems
 - Properties of dependable systems
 - Redundancy, process resilience (see T8.2)
 - Reliable RPCs (see T8.3)
- Fault Tolerance Backward recovery using checkpoints.
 - Tradeoff: independent vs coordinated checkpointing
- Fault Tolerance Recovery using Write-Ahead-Logging
 - Balances the overhead of checkpointing and ability to recover to a consistent state

Additional Material in the Book

- Process Resilience (when processes fail) T8.2
 - Have multiple processes (redundancy)
 - Group them (flat, hierarchically), voting
- Reliable RPCs (communication failures) T8.3
 - Several cases to consider (lost reply, client crash, ...)
 - Several potential solutions for each case
- Distributed Commit Protocols T8.5
 - Perform operations by all group members, or not at all
 - 2 phase commit, ... (covered in last lecture)
- Logging and recovery T8.6