

Distributed Systems

Christian J. Rudder

January 2025

Contents

Contents	1
1 Introduction	5
1.1 Transactions and Concurrency Control	6
1.1.1 Optimistic Concurrency Control (OCC)	6
1.1.2 Two-Phase Commit (2PC)	10
1.1.3 Three-Phase Commit (3PC)	12
Bibliography	14

This page is left intentionally blank.

Big thanks to **Professor Ioannis Liagouris**
and **Dr. Anna Arpaci-Dusseau** for teaching CS351: Distributed Systems
at Boston University [\[1\]](#).

All illustration contain original assets.

*Disclaimer: These notes are my personal understanding and interpretation of the course material.
They are not officially endorsed by the instructor or the university. Please use them as a
supplementary resource and refer to the official course materials for accurate information.*

Prerequisites

This text assumes the reader has a basic understanding of computer science and programming. It will also assume they are somewhat familiar with computer architecture and operating systems at a high level. The text will review these concepts briefly for completeness, but it will not try to teach them from scratch or provide a full understanding of these topics.

The main focus will be on distributed systems, and will touch on:

- **Concurrency and Parallelism**
 - Concurrency, Parallelism, Threads
- **Consistency and Fault Tolerance**
 - Consistency, Fault-tolerance, Atomicity
- **Distributed Systems and Coordination**
 - Asynchrony, Coordination, Logical Time, Snapshots
- **Consensus Algorithms**
 - Raft, Paxos, Consensus
- **Replication and Data Management**
 - Replication, Sharding, Cluster
- **Protocols and Computing Models**
 - RPC, 2PC, Broadcast
- **Technologies and Tools**
 - MapReduce, Spanner, Dynamo, GFS, TLA+, Golang

— 1 —

Introduction

1.1 Transactions and Concurrency Control

1.1.1 Optimistic Concurrency Control (OCC)

Say the backbone of our stock trading application is a distributed database. The system may conduct complicated stock trades based on server stock prices.

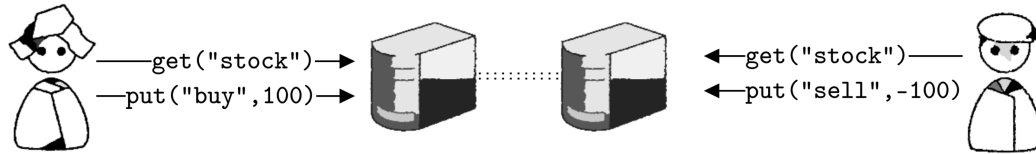


Figure 1.1: Two users checking the stock prices and making a trade based on such information.

It is critical that the stock price is consistent through all servers, and more important that if any trade fails, the system can recover to a consistent state.

Definition 1.1: Transaction

A **transaction** is a sequence of operations that are treated as a single unit of work. A transaction must satisfy the **ACID** properties:

- **Atomicity:** A transaction is either fully completed or dropped entirely.
- **Consistency:** A transaction must leave the database in a consistent state.
- **Isolation:** Transactions must be isolated from each other.
- **Durability:** Once a transaction is committed, it remains so even after system failure.

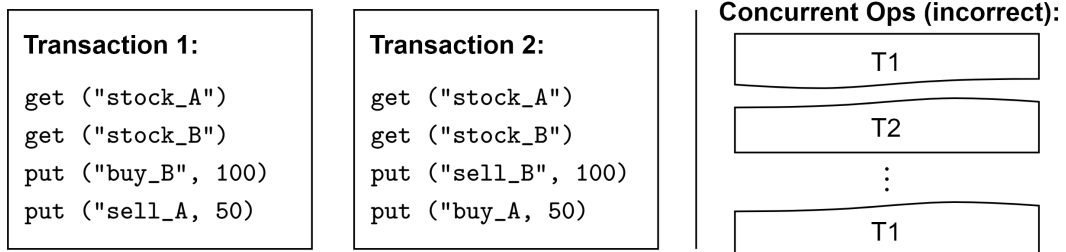


Figure 1.2: Two transactions which whose operations are interleaved.

Interleaving transactions **violates the isolation property**. This is problematic in Figure (1.2) as T2's (transaction 2) operations may depend on the server state before T1's transaction. Additionally partially completed transactions leave the system in an inconsistent state, violating **atomicity** (e.g., T1's "buy_B" fails, but "sell_A" succeeds).

We discuss another consistency model which will help us in this settings:

Definition 1.2: Serializability

Serializability is a strong consistency model that ensures the outcome of concurrent transactions is the same as if they were executed in some sequential (serial) order.

This differs from **linearizability**, which focuses on the real-time ordering of individual operations. Serializability instead concerns the logical order of **entire transactions**, independent of their timing.

However, **strict-serializability** does care about real-time ordering in addition to the logical order of transactions. This implies linearizability, but not vice versa.

We consider the following model to help us preform transactions:

Definition 1.3: Optimistic Concurrency Control (OCC)

Optimistic Concurrency Control (OCC) assumes conflicts are rare and proceeds without locking. It follows four main steps:

- **Prepare:** The system reads the transaction request and creates a backup or temporary copy of the state.
- **Modify/Validate:** The transaction modifies the temporary state. Then The system checks whether the transaction is **serializable**.
- **Commit/Rollback:** If valid, commit; Otherwise, abort transaction and rollback to previous state.

This only solves **isolation**, as it does **not** guarantee atomicity.

Definition 1.4: Transaction Coordinator and Database Servers

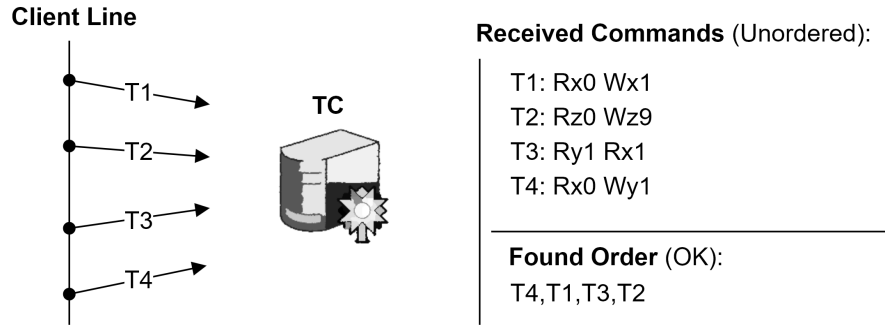
OCC maintains two necessary components:

- **Transaction Coordinator (TC):** The validation server responsible for checking whether a transaction is **serializable**. It receives requests from clients and responds with either:
 - OK: if the transaction is serializable,
 - ABORT: if it conflicts with prior transactions.
- **Database Servers (DB):** Receives transaction operations, executing it on local state. Then, on **OK**—commit state, on **ABORT**—rollback to the previous state.

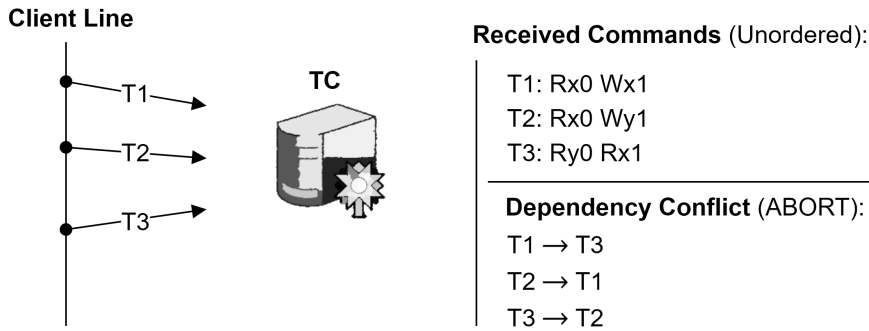
We consider one model, which where multiple clients interact with one TC.

Example 1.1: Centralized OCC

Consider these two examples with a single TC and multiple clients on the network line:



Here, clients on the line send transactions T1–T4 to the TC. The TC then checks the transactions for serializability. In this case an order is found (T4,T1,T3,T2), which the TC communicates to the DBs to commit.



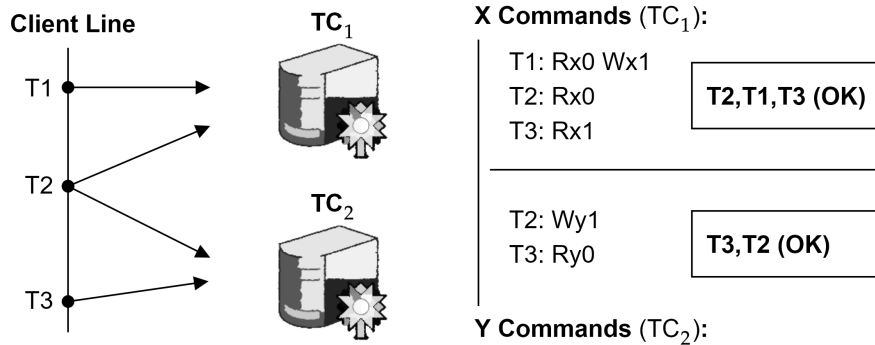
Here, transaction requests, T1, T2, and T3, do not have a serial order. As we build, T1 → T3 (T1 then T3) makes logical sense. Then T2 → T3, giving us the order T2 → T1 → T3; However, it appears T3 must come before T2. This causes a cycle (T2 → T1 → T3 → T2). Hence, the TC must must abort all transactions. ■

Tip: If familiar with **Directed Acyclic Graphs (DAG)**, one can think of the transactions as nodes and the edges as the dependencies between them. If the graph is a DAG, then there is some serial order (OK). If not, then there is a cycle, so we must abort.

Though we run into an issue when there are multiple TCs.

Example 1.2: Distributed OCC

Consider two TCs responsible for different parts of our system data.

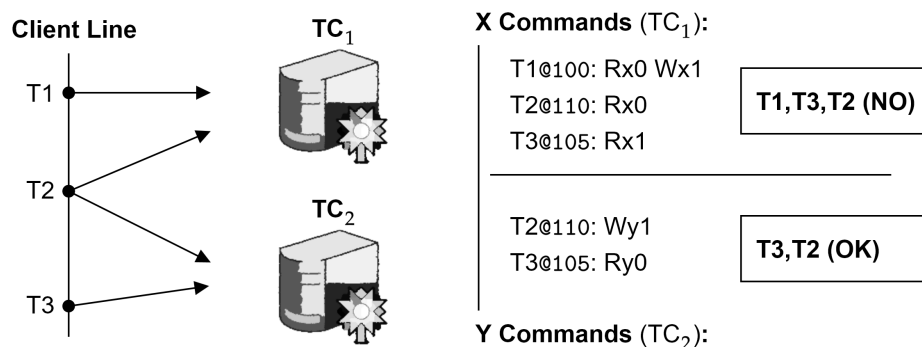


The problem occurs as TC₁ and TC₂ pick **different serial orders** for the transactions. ■

Theorem 1.1: Timestamping Distributed OCC

Timestamping is a method where each transaction is assigned a unique timestamp (ID), which aids order agreement between TCs. The downside: **unnecessary aborts**.

Example 1.3: Timestamping Distributed OCC



Timestamps (@#) are only serve as **IDs**. Here TC₂ OKs the order (T3, T2). TC₁ sees this, enforces the order, but is not able to serialize (T1, T3, T2). Hence, it aborts (NO). ■

1.1.2 Two-Phase Commit (2PC)

We introduce another method to help us with this problem, though it **does not ensure isolation**:

Definition 1.5: Two-Phase Commit (2PC)

Two-Phase Commit (2PC) ensures **atomicity**. The client sends the transactions to the DBs (participants). There after, the client tells the TC start the commit process, involving two phases:

- **Prepare Phase:** The coordinator sends a prepare request to all DBs; Each respond:
 - YES: if it can commit the transaction.
 - NO: if it cannot commit the transaction.
- **Commit Phase:** If all voted YES, the coordinator sends a **COMMIT** request to all DBs. If any participant voted NO, the coordinator sends an **ABORT**. The TC then responds to the client with the final outcome of the transaction.

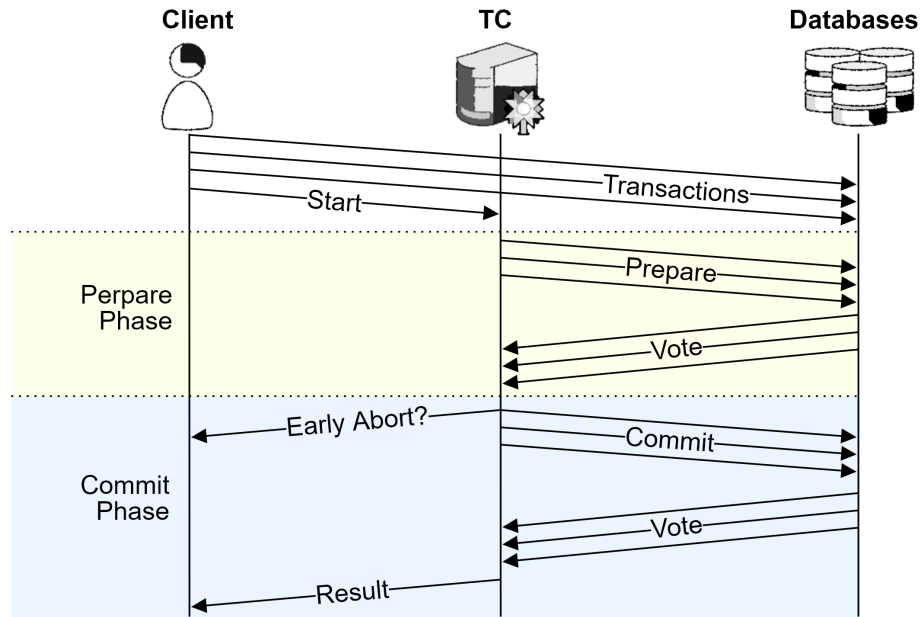


Figure 1.3: Two-Phase Commit (2PC) process. The client sets up the transaction with the TC and DBs. Then, the TC starts the commit process, starting with the prepare phase, and then ends with the commit phase.

Though there is a pitfall with this method:

Definition 1.6: 2PC Blockout

In most cases if there's a timeout the TC or DBs will abort the transaction. However, if the TC times out after a DB has voted **YES**, the DB is left in an uncertain state, and must **wait** for the TC to respond.

Preparation Phase	
Case	Action
TC timeout for yes/no vote	Abort — transaction coordinator did not receive votes in time.
DB timeout for prepare	Abort — database did not receive prepare request in time.
Commit Phase	
Case	Action
DB timeout for commit/abort and DB voted NO	Abort — since DB already voted NO, it's safe to abort.
DB timeout for commit/abort and DB voted YES	Block — DB must wait for TC decision to preserve atomicity.

Definition 1.7: 2PC Persistent & Volatile State

In the Two-Phase Commit (2PC) protocol, both the Transaction Coordinator (TC) and Database (DB) participants must persist critical information to recover correctly after a crash.

- **Database (DB):**
 - Must persist the result of any vote (YES or NO).
 - If the DB voted YES and then crashes, upon recovery it must contact the TC to learn the final decision (COMMIT or ABORT).
- **Transaction Coordinator (TC):**
 - Must persist the result of any vote and the final decision (COMMIT or ABORT).
 - If the TC crashes, it must resume the commit protocol from where it left off.

1.1.3 Three-Phase Commit (3PC)

2PC's main problem was **availability**. We can fix this by adding an additional phase:

Definition 1.8: Three-Phase Commit (3PC)

Three-Phase Commit (3PC) is an extension of 2PC that adds a third phase to avoid blocking in case of failures:

- **CanCommit Phase:** The coordinator sends a **CANCOMMIT** request to all participants. Each participant responds with either **YES** or **NO**.
- **PreCommit Phase:** If all participants respond with **YES**, the coordinator sends a **PRECOMMIT** request to all participants. Participants prepare to commit sending back an acknowledgment (**ACK**).
- **DoCommit Phase:** If all participants **ACK** the precommit, the coordinator sends a **COMMIT** request. Otherwise, it sends an **ABORT** request.

Theorem 1.2: 3PC Non-Blocking Nature & Self-Resolution

The reason this fixes the blocking issue in 2PC lies within the **PRECOMMIT** phase. In regular 2PC, if a DB votes **YES** it must wait for the TC to respond as it is uncertain whether another DB has committed or not.

In 3PC, such DB need not worry, as no other DB can commit until the TC sends a **PRECOMMIT** request. Henceforth, DBs can safely either terminate or resolve between themselves that if everyone else has the **PRECOMMIT** request, they can commit.

Theorem 1.3: 3PC Persistent & Volatile State

Both the TC and DBs may crash at any time, for which they must persist the state of the transactions. Though it is fully possible for the TC to only persist the original transaction request, and instead query all DBs for their state to recover the transaction.

Theorem 1.4: 3PC Tradeoffs

- (−) If the network becomes partitioned, inconsistencies may arise within self-resolution.
- (−) Adds another round trip to the commit process (3 in total).
- (+) Removes the blocking issue of 2PC, as the TC can always recover from a crash.

Below is a diagram of the 3PC process:

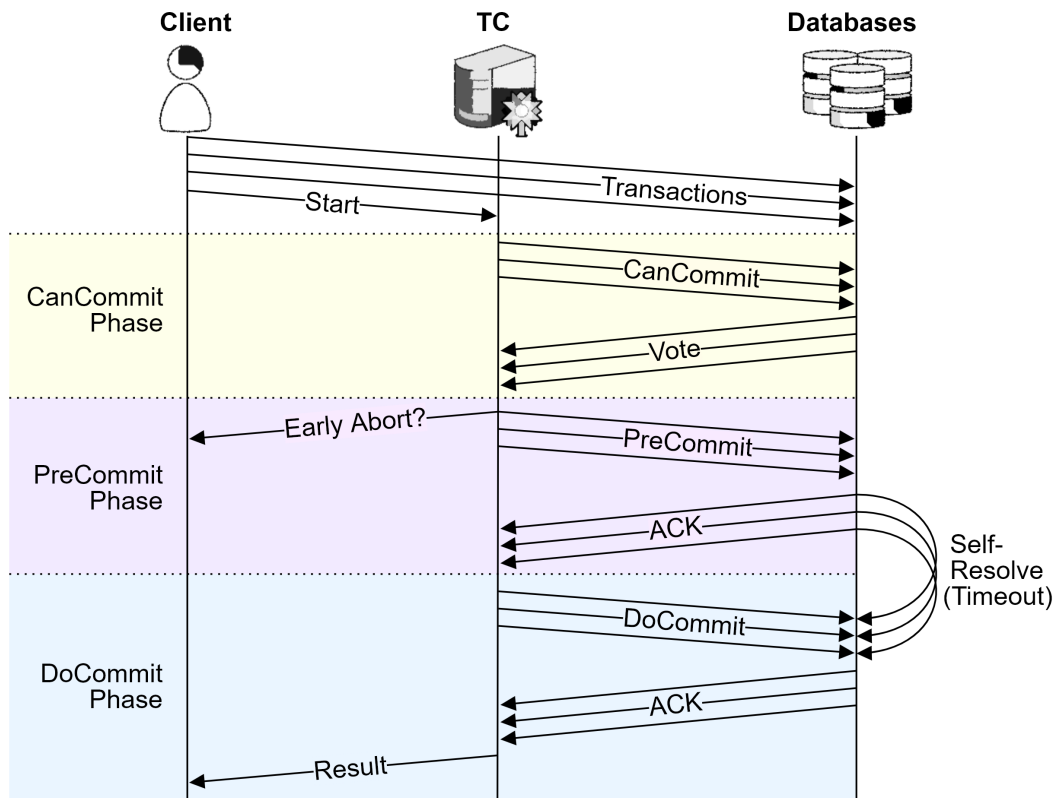


Figure 1.4: Three-Phase Commit (3PC) protocol message flow. The transaction progresses through three phases: **CanCommit**, where participants vote; **PreCommit**, where participants acknowledge readiness to commit; and **DoCommit**, where the final decision is executed. In the event of coordinator failure, participants may invoke **Self-Resolve** after a timeout, based on whether they received a **PreCommit** message, ensuring non-blocking recovery. **Note:** The Databases do not communicate with the client after resolution, though this could depend on the implementation.

Bibliography

- [1] Ioannis Liagouris. Cs351: Distributed systems. Lecture notes, Boston University, Spring Semester, 2025. Boston University, CS Department.