

Distributed Inventory System

Project Distributed Systems

for the

Master of Science

from the Course of Studies Computer Science

at the University of Stuttgart

by

Luca Schwarz, Johannes Hartmann, Patrick Baisch

11.02.2024

Due Date

16.10.2023-11.02.2024

Student ID

3731224, 3743267, 3125018

Reviewer

Marco Aiello

Contents

Acronyms	II
List of Figures	III
1 Introduction	1
2 Project Requirement Analysis	2
2.1 Architectural Description	2
2.2 Dynamic Discovery of Hosts	3
2.3 Fault Tolerance	3
2.4 Leader Election	4
2.5 Ordered Reliable Multicast	4
3 Component Description	6
3.1 Node	6
3.2 Leader	6
3.3 Client	7
4 Discussion	8
4.1 Crash	8
4.2 Fail Stop	8
4.3 Byzantine	9
4.4 Voting	10
4.5 Ordered Reliable Multicast	11
5 Summary	12
Appendix	13

Acronyms

BFT	Byzantine Fault Tolerance
UUID	Universally Unique Identifier
FIFO	First In First Out

List of Figures

2.1	Architecture Diagram	2
-----	--------------------------------	---

1 Introduction

As part of the course project an inventory system shall be developed. The system will be a client server architecture where the servers are responsible to keep the data and process updates regarding goods and stock information. The clients can request the currently available goods and the amount and send update information if new goods are available or if goods are taken out of stock. For the updates it's important to ensure strong consistency and ordering of events such that all clients have the current information. Within the servers it shall be possible to add additional nodes to the system dynamically and to handle different failure cases while still being able to process requests.

2 Project Requirement Analysis

2.1 Architectural Description

Client-Server Architecture: Our system will adopt a client-server model as seen in figure 2.1, which consists of multiple clients interfacing with a cluster of server nodes. The clients are responsible for querying current stock levels and submitting updates for processing. The server nodes serve these stock level requests and manage inventory updates. Amongst the server nodes, a leader will be elected to coordinate updates and ensure consistency across the distributed system. Our main focus is the strong consistency, therefore if a new request is past from the client to a server, first the server reaches out to the leader, which will commit the change and pass it to all other servers. Writing requests are always passed to the leader, before they take place. Reading requests can be handled from a server node, the leader is not involved in this actions. When the leader dies, no bussiness logic can happen, except for reading, before a new leader is chosen.

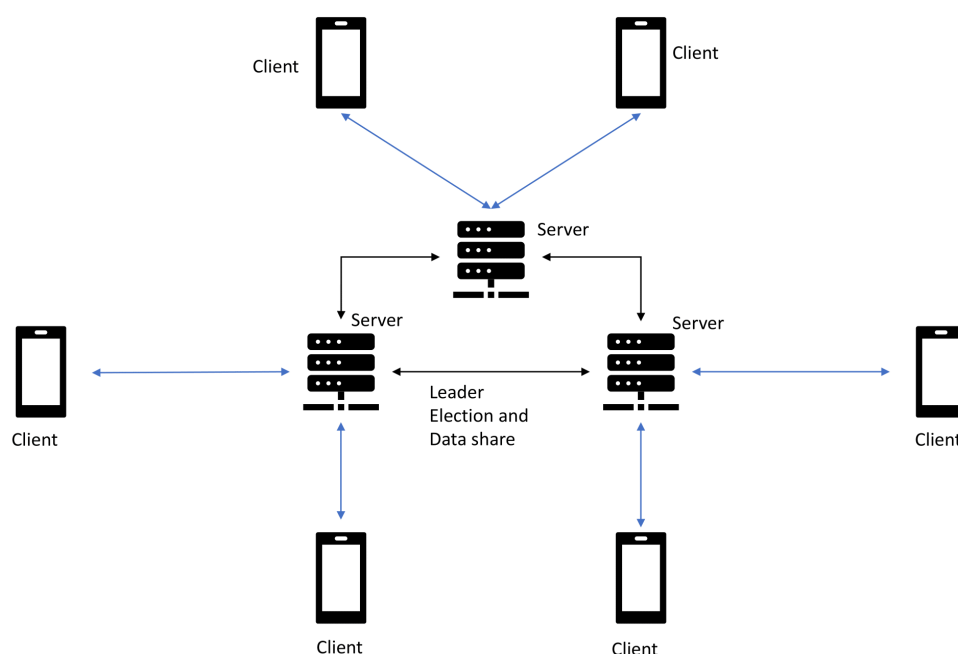


Figure 2.1: Architecture Diagram

2.2 Dynamic Discovery of Hosts

Server-side: Upon initiation, each server node will broadcast its presence and listen for existing members of the system to construct a current view of the cluster. The broadcast message from the joining server is only responded by the leader of the system. The message contains an array with all ip addresses currently active within the system, as well as the information which node is the leader.

When the servers are started, no leader is chosen at this point of time. In this case, the first node which does not get a reply of his broadcast declares himself as leader and answers the requests of the other nodes. To ensure, that only one leader is chosen at the starting point, a random sleep is implemented, which ensures the existence of only one leader.

This dynamic discovery protocol allows the system to scale horizontally without manual configuration.

Client-side: Clients are designed to automatically detect server nodes in the system. As well as the server-side implementation a broadcast message is sent, but not with the intention to join the cluster. In the Client-side case the ip address is not added to the server node pool. The reply to the broadcast is done by the leader, sending an array of all members to the client. The client uses one randomly chosen address to communicate with. If the reply of an request by an client is not answered by the node, another discovery will be done to get an updated version of the list with nodes.

This enables seamless interaction with the inventory system, ensuring that clients can always locate a server node to process their requests.

2.3 Fault Tolerance

Our system is engineered to handle different types of failures, ensuring continuous operation:

Leader Failure: If the current leader server fails, the remaining servers will initiate a leader election to select a new leader. In the meantime no operation is permitted, because our system uses strong consistency. As soon as a new leader is elected, the system returns to normal process. This ensures that the system become available after a crash of the leader.

Server Crash: In the event of a server crash, client requests are automatically redirected to other operational server nodes. A health check mechanism and a server list update protocol ensure that clients and servers are aware of the available nodes in real-time. **Luca** When a

node does not send a Heartbeat, it will be removed from the list by every server, including the leader, so the normal processing can continue.

2.4 Leader Election

Consensus on Updates: Write operations, critical for inventory synchronisation, are managed by the leader node. Also new nodes which want to join the group exchange information with the leader. The election Algorithm will be the bully Algorithm. When an election is triggered (the heartbeat of the leader is missing), the first Server, which realises it, starts an election, using an ID as information multicastet to every other server. If another Server has an higher id, he claims to become the leader, “bullying” the server with a lower ID and denying his election. As ID we use the IP Address as string to compare. This will continue untl the server with the highest ID wins the election. He then sends a message, declaring himself as new leader and the workflow can continue.

2.5 Ordered Reliable Multicast

Our system uses total ordered reliable multicast to ensure inventory updates, like purchases or sales, are processed sequentially and data consistency is maintained over all nodes. This also prevents that a host has outdated information.

In our case the leader acts as the sequencer. All business messages are sent to the leader and the leader then sends the messages to all hosts via udp multicast. Each message contains a sequence number and the actual content. When a host receives a message from the leader the message is put into a log and an acknowledge message is sent back to the leader. Upon receiveing an acknowledge message the leader adds the sender to the acknowledge list and checks if all hosts acknowledged the message. It's required that all hosts acknowledge a message as we strive for full consistency over all active hosts. After receiveing enough acknowledgements the leader sends a commit message to all hosts triggering that the message is processed.

For each received commit message the host checks if the sequence number is strictly one higher than the sequence number kept by the host. If that's the case the message is processed. If that's not the case the host missed messages. To ensure consistency and ordering of the messages it requests from the leader all messages starting from it's current sequence number up to the received sequence number. The leader then sends the

messages in the order of the sequence numbers. The messages can then be consumed by the host.

3 Component Description

3.1 Node

Each node in the system will check on startup if there are other nodes already running within the network. This is done via a broadcast message containing the nodes ID and IP Address. If after a timeout no host responded the system assumes it's the first node and becomes leader. If other systems are already active the leader responds to the broadcast with the current list of hosts, it's own IP (the leader IP). The new host updates its own host list with the information received from the leader and sets the leader IP address.

Each node sends regularly heartbeats to all other nodes. It then checks if there are missed heartbeats from other nodes. If for a node the heartbeats are missed for too long it removes the host from the hosts list. If a node detects a failure of the leader a new election is started.

Each node listens for business requests from clients. If a business request is received by a host not being leader it forwards the request to the leader and waits for a response. Which is then forwarded to the requesting client.

3.2 Leader

The first leader is determined by the first host in the system. In case the leader fails a new leader is elected out of all nodes using the bully algorithm to determine a unique leader.

The leader behaves the same as all other nodes regarding heartbeat and processing of committed messages.

In addition to the above the leader responds to broadcast messages from new hosts and clients with the current hosts list.

Processes all business logic messages and distributes them using ordered reliable multicast to the other nodes. After a message is committed and processed the leader responds to the host which sent the message.

3.3 Client

The client is a web interface that can be used to interact with the warehouse system. Since only a handful of functions are needed, the interface was created using Flask, a micro web framework written in Python. Furthermore Bootstrap was used to create a simple but reliable user interface without the need of custom CSS. In addition, the interaction between the two components is simple and straightforward to implement.

The interface offers four functions on three pages: All available items are listed, an item can be created, and items can be bought or sold. Triggering a function causes the client to send a request to all currently active hosts. An active host is selected at random. The requested information is then sent to this host and the response is awaited.

The requests consist of a type that triggers which server routine is requested, followed by the corresponding metadata, which comes from the user via a form. The forms already have validators so that the data does not need to be checked again.

4 Discussion

4.1 Crash

As mentioned in section [2.3](#) our system is able to handle the crash of a server node, as well as the crash of the leader node. The procedure of the crashes differs, whether a node dies, or the leader dies. In case of a node dying, the system can operate as if nothing happens. Internally the other servers realize, that there is no heartbeat of the dead node, After 2 seconds without a heartbeat, the server node is taken out from the cluster and in the meantime the system operates as usual. The benefit of this approach is, that there is no downtime at all. But it can happen, that a request from a client is not registered and it has to be requested again, in order to handle the change.

A slightly different approach takes place, when the leader dies. In that case the system is not operational, until a new leader is elected. As soon as a node realizes the death of the leader, the election will be started and a new leader will be chosen. This ensures strong consistency, even without the leader but at the cost of a little downtime, where the requests from the clients will be paused, until the new leader is operational. Requests that are already placed but not finished, will be lost in order to keep the system consistent, therefore the user has to place the requests again.

4.2 Fail Stop

Fail Stop is the stop of a system, if a node/component is not operational. One of the primary advantages of fail-stop lies in its simplicity and predictability. When a component or node within a distributed system fails, it halts all operations immediately and unequivocally. This clear-cut behavior simplifies fault detection and recovery processes, as the failed component can be swiftly identified and isolated without ambiguity. Consequently, fail-stop facilitates efficient fault management strategies, minimizing downtime and enabling rapid system restoration.

However, fail-stop is not without its limitations and disadvantages. One notable drawback is the potential for abrupt service disruptions and loss of progress. Since fail-stop entails an immediate cessation of operations upon failure, ongoing tasks or transactions may be

abruptly terminated, leading to potential data loss or service interruptions. This characteristic can be disruptive, particularly in scenarios where the system's continuity and uninterrupted operation are paramount.

In order to fulfill our strong consistency, we did not implement a dedicated fail stop procedure. But the server nodes will halt its operations if the leader is dead to prevent the creation of inconsistent data.

4.3 Byzantine

Byzantine Fault Tolerance (BFT) protects the system from malicious attacks. Derived from the Byzantine Generals' Problem, which illustrates the challenge of achieving consensus among mutually distrustful parties, BFT mechanisms fortify distributed systems by enabling them to operate seamlessly even in the presence of faulty or malicious nodes. The Duality of Byzantine Fault Tolerance: A Double-Edged Sword

Byzantine Fault Tolerance (BFT) embodies a dual nature in the realm of distributed systems, presenting both positive and negative aspects to its implementation. This essay explores the dichotomy of BFT, highlighting its benefits and drawbacks in ensuring the reliability and integrity of distributed systems.

On the positive side, the implementation of BFT brings forth a robust defense mechanism against malicious attacks and system failures. BFT algorithms enable distributed nodes to coordinate and reach consensus, even in the presence of Byzantine faults. This capability ensures the integrity and consistency of the system, enhancing its resilience against adversarial conditions.

Moreover, BFT offers fault isolation and containment mechanisms, minimizing the impact of faulty or malicious nodes on the overall system. By detecting and isolating Byzantine faults, BFT prevents the spread of disruptive influences, thus maintaining the operational continuity of distributed systems. This aspect of BFT is crucial for ensuring system reliability and availability, especially in critical applications where uninterrupted operation is paramount.

However, the implementation of BFT also comes with inherent challenges and drawbacks. One notable negative aspect is the increased communication overhead associated with consensus protocols. BFT algorithms typically require multiple rounds of message exchanges among nodes to achieve agreement, resulting in elevated network traffic and latency. This overhead can impact the performance and scalability of distributed systems, particularly in large-scale environments or under high network load conditions.

Furthermore, achieving Byzantine Fault Tolerance often involves a trade-off between fault tolerance and system efficiency. While BFT mechanisms excel in ensuring fault tolerance and resilience, they may introduce complexities and resource overheads that can hinder system performance and scalability.

In our case we decided, that we do not any sort of BFT to have higher speed. The procedure we use to ensure strong consistency is already slowing down the system. In production and widely used system, implementing BFT is crucial concerning the security.

4.4 Voting

In order to select a leader as described in 2.1, a voting strategy is crucial. We chose the bully Algorithm to vote for a leader. In this algorithm, any node can start the election, but the node with the highest ID, whatever ID is chosen (e.g. Universally Unique Identifier (UUID) or IP Address), will win the election by “bullying” the nodes with a lower ID. If the comparison with any other ID is greater instead of greater or equal, if two nodes have the same ID, the node which voted first will win.

The Bully Algorithm provides a decentralized approach to electing a leader in a distributed system. By allowing nodes to autonomously determine the leader without relying on a central authority, the Bully Algorithm enhances system resilience and scalability. This decentralized nature ensures that the system remains operational even if certain nodes fail or become unreachable, thereby improving fault tolerance and system robustness.

However, one limitation of the Bully Algorithm is its susceptibility to network partitioning or communication failures. In scenarios where network partitions occur, nodes may incorrectly elect multiple leaders within different partitions, leading to inconsistencies and conflicts within the system. Additionally, if communication failures prevent nodes from exchanging messages during the leader election process, the algorithm may result in delays or failures to elect a leader, compromising the system’s responsiveness and reliability.

We chose the Bully Algorithm for its decentralized approach to leader election, enhancing fault tolerance and scalability in distributed systems. As a bigger System is developped, an other algorithm should be implemented due to factors such as scalability limitations, susceptibility to network partitions, and the need for more complex fault tolerance mechanisms tailored to their specific requirements and workload characteristics.

4.5 Ordered Reliable Multicast

Ordered reliable multicast provides a method to distribute information over a known set of nodes in a reliable manner and ensures ordering of the messages depending on the requirements. If no ordering mechanism is implemented the message order might differ on the hosts for example caused by network delays.

An alternative order is First In First Out (FIFO) in this case the order of the sender of the messages has to be ensured over all nodes. This can be achieved using vector clocks to track the sequence of messages for each node. Other order schemes are happened before relationships where it can be told if a message happened before an other message if they have a causal relationship but nothing can be said about messages which happened independent of each other.

With Total order all messages over all nodes have the same order. To achieve this a central system is needed which generates the sequence e.g. by computing sequence numbers.

As already noted above, we selected total order for our implementation and defined the leader to provide a monotone sequence number. In our implementation all messages being it actual messages with content or acknowledge / commit messages are sent to the multicast group. This leads to high network traffic with growing nodes. An improvement could be that acknowledge messages are only sent to the leader over unicast as it's not of interest for the other nodes. Also, if a node detects that it missed messages then again all messages are sent over the multicast group even if they only concern the node which missed them.

5 Summary

Overall we built a warehouse system consisting of a clients and servers. The server nodes can scale horizontally and find other nodes via dynamic discovery. To ensure consistency over all nodes we implemented communication over ordered reliable multicast with total order. The server nodes elect a leader which acts as central point for onboarding new nodes and the ordered reliable multicast. We implemented the bully algorithm for the election. Regarding fault tolerance we implemented mechanisms to detect a system Crash and handle it in a way that does not lead to inconsistencies. An algorithm to handle [BFT](#) is not implemented in our system.

To sum everything up we built a distributed system implementing a minimum set of protocols and features ensuring scalability and stability. To show the functionality we implemented a small business application on top. In order to use the system in a productive manner it would be necessary to implement further fault tolerance mechanisms and enhance some protocols regarding time and number of messages exchanged.

Appendix

words: 3109

Github: <https://github.com/jojo221119/distributedInventoryManagement>