

## Application Level Scalable Leader Selection Algorithm for Distributed Systems

Zahir Sayyed\*

Software Engineer, Jamesburg, New Jersey, USA.

\* Corresponding Author Email: [sayyedzahir1@gmail.com](mailto:sayyedzahir1@gmail.com) ORCID: 0009-0004-6555-3228

### Article Info:

DOI: 10.22399/ijcesen.3856

Received : 29 July 2025

Accepted : 06 September 2025

### Keywords

Leader Election Algorithm,  
Distributed Consensus,  
Fault Tolerance,  
Scalability in Distributed Systems,  
Consensus Protocols,  
Failure Recovery Mechanisms,

### Abstract:

Leader selection is a critical issue in distributed systems where a single node is elected from a number of nodes to perform specific tasks or to provide system consistency. This article proposes the design and implementation of a scalable leader selection algorithm that ensures efficiency in large-scale distributed environments. This algorithm is designed to minimize communication overhead, tolerate node failures and scale with the size of the system. Leader selection is the most important issue for distributed systems since it guarantees the coordination and efficiency of the system. In this paper, we propose a scalable leader selection algorithm that can adapt to the system size and network condition changes. We will discuss its design, implementation challenges and performance..

## 1. Introduction

Distributed systems are made up of multiple parts that work together to achieve a common goal. One of the challenges is ensuring that a single node server act as a leader. A leader can be any node that is responsible for tasks at the system level. For example, it is responsible for the task scheduling, the resource management, and the fault tolerance. The leader is selected using leader election algorithms. The leader election algorithms that are available today have a disadvantage, which is they are not scalable, they have a high communication overhead, and they are not fault tolerant. The goal of this paper is to provide a way to make a leader election algorithm the most scalable and with the minimal overhead.

Leader election is a powerful tool for improving efficiency, reducing coordination, simplifying architectures, and reducing operation (1). Leader selection algorithms may encounter scalability problems in a larger distributed system. These problems may include communication overhead and other issues that would need to be addressed. In this paper, a scalable leader election algorithm is presented, which is designed to minimize resource usage and communication overhead while maintaining fault tolerance and reliability.

Distributed systems are the fundamental building block of the current computing landscape, from cloud services to financial systems to

telecommunication networks to content delivery networks. At their very core lies a challenge of determining and maintaining leadership over a group of systems. If a system lacks a mechanism for finding a leader, it will have poor data consistency, request routing, failure detection, and recovery coordination.

In the current world of distributed systems, many systems depend on external coordination services, such as ZooKeeper, etcd, or Consul, for managing leader election. However, the biggest problem with external coordination services is the additional baggage it adds by means of an extra system to maintain. This can lead to many new types of failures and even performance bottlenecks that limit scalability.

The research gap being addressed is what happens when an external coordination service is not available, or not suitable, or even serves as the slowest link in the system. With the help of the proposed novel Algorithm for event-driven Application-Level Scalable Leader Selection Algorithm for Distributed Systems, which integrates leadership coordination to the application level, thereby overcoming the need for external coordination while providing the highest level of flexibility and superior performance and return.

Consider the following scenarios where traditional coordination approaches fail: A manufacturing firm uses thousands of edge devices in factories with

unreliable network connections. External coordination services require continuous connection to central servers, rendering them unsuitable for this application. Financial institutions must obey stringent data sovereignty rules and cannot utilize third-party coordination services that might send traffic between jurisdictions. Furthermore, IoT deployments and embedded systems cannot maintain connections to heavyweight coordination services while performing their primary objectives. In each case, our application-level approach solves the issue by incorporating leader election into the application's communication patterns, necessitating only a small number of additional resources to guarantee strong leadership guarantees.

Application-Level Scalable Leader Selection Algorithm is an essential innovation in distributed systems coordination. This approach introduces the concept of leadership directly to the application layer significantly improving performance, resiliency, and flexibility and reducing dependencies. Its broad applicability across numerous industries and technical fields showcases the breakthrough it represents. The innovation addresses crucial challenges of the current state of distributed computing, especially as systems move beyond conventional data centers to edge environments, resource-limiting devices, and highly regulated settings. The system's capability to offer solid leader selection in tough situations makes it possible to develop completely new classes of distributed applications previously deemed too uncertain to create.

## 2. Related Work

The election of leader is based on the lease involves the issuance and extension of a lease taken by a particular node that has the right to become a leader. In terms of the lease, the node is the leader of the other nodes. If a valid lease is not extended in time, the other node will be able to become the leader when it has a valid lease. Leaders' lease does not give other nodes the right to become leaders. It is only possible that the lease of the node was not renewed because of the crash or communication problems between nodes. Thus, there is only one leader due to the availability of a valid lease.

The kind of leadership that is shard-based is utilized in systems that separate data into distinct units or shards. Different leaders manage and maintain data from separate shards, ensuring parallel processing. Distributed processing is both scalable and fault-tolerant because each shard has its own independent leader. For example, AWS DynamoDB and EBS use

this strategy to manage data across a cluster efficiently. Coordinating between shards and maintaining consistency can be difficult if the shards are split, merged, or if there are changes in leadership.

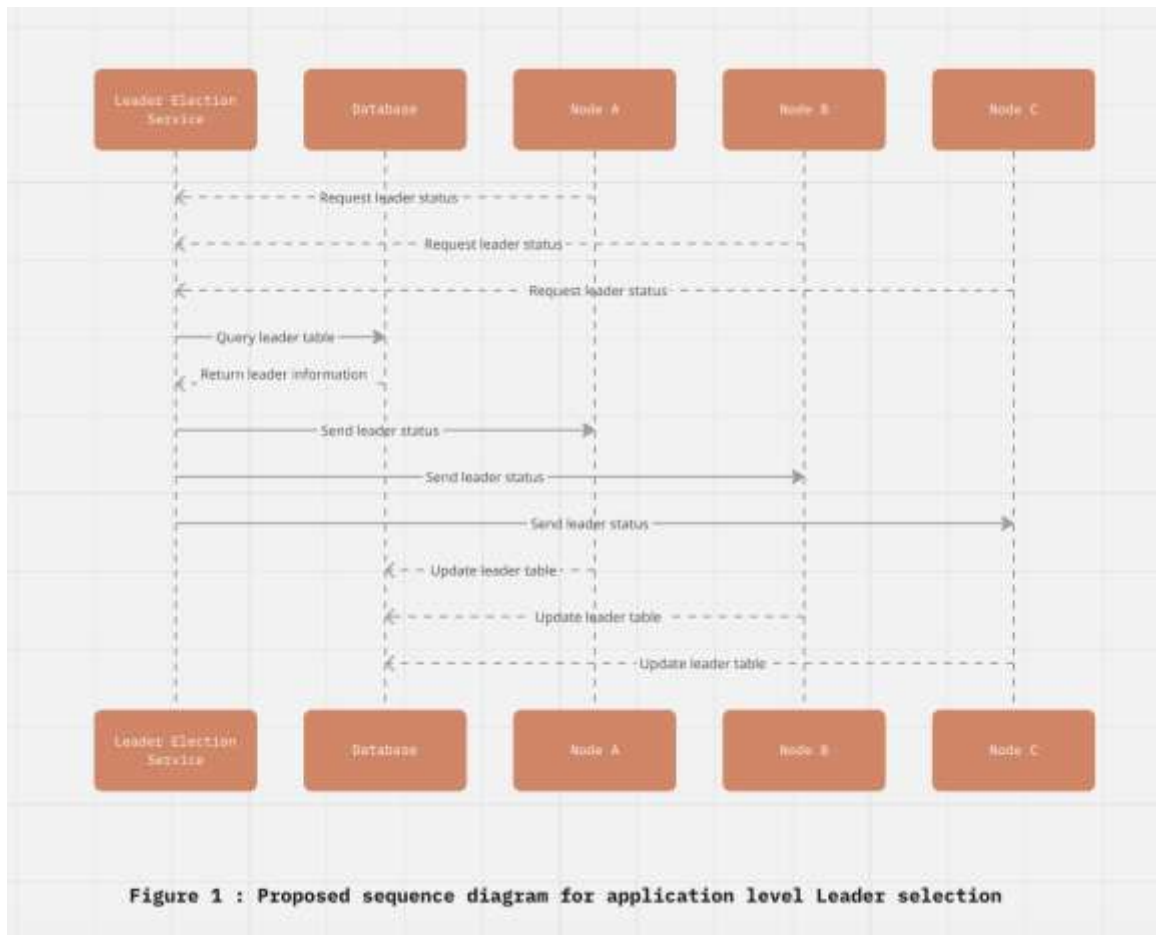
Consensus algorithms like Raft and Paxos are used to elect a leader fault-tolerantly by reaching an agreement among a quorum of nodes. These mechanisms guarantee that only one node becomes the leader and all others agree on the designation even if failures or partitions occur. This method is great for systems that require strong consistency and reliability, as it helps avoid conflicting decisions of leadership. However, it may add communication overhead since it is necessary for an agreement to be reached.

Without an effective leader election mechanism, distributed systems may face significant issues regarding data inconsistency and conflict resolution. Imagine a scenario where two nodes believe they are the leader, leading to data divergence and potential conflicts that are hard to reconcile. Leader elections ensure that only one node becomes the leader, thereby maintaining harmony and consistency across the system (3).

External services like ZooKeeper and etcd offer advanced features for leader election, including ephemeral nodes and distributed locks. Consequently, a system can delegate the election of a leader to these external services to reduce the complexity of its functioning. The systems oversee the maintenance of a single source of truth and also offer automatic failover in case the leader fails. The services are dependable and in use in diverse environments, but they also demand operational overhead, result in latency to make decisions on leaders and pose a challenge for the engineers to manage properly to prevent them from becoming a bottleneck or a single point of failure.

## 3. Proposed Approach

Application-level leader election is a strategy where the application logic itself is responsible for selecting a leader between distributed nodes as opposed to using external coordination services or frameworks such as Akka Cluster Singleton. This method is particularly useful in systems where simplicity, reduced dependencies, and tight integration with application-specific requirements are critical. Leader election ensures that all nodes in a distributed system agree on a single node to act as a coordinator or leader/commander to manage tasks and communication (4).



**Figure 1** shows the overall leader selection sequence, One typical approach to applying application-level leader election is to use distributed mutexes or locks. Each node in the network makes an attempt to acquire the lock of a shared resource; when one of them successfully acquires the lock, it acts as the leader. As a result, only one node can be the leader at any given time, thereby leading to a consistent and well-synchronized network. In case the node that is functioning as the leader unexpectedly crashes or releases the lock, some other node may compete and acquire the lock to be the new leader. This feature makes the network fault-tolerant and highly available.

In comparison to other leader election approaches, the application-level leader election method provides more flexibility and control to developers. The technique gives power to developers to modify the election process according to their needs and requirements of their specific application. Moreover, this method reduces the need for external coordination services, which are required in traditional methods, and naturally, this method is more relaxed to implement where these services are not needed or desired. On the contrary, this technique requires the developer's attention to handle specific scenarios, such as network partitions and the failure of a node. Avoiding scenarios such as

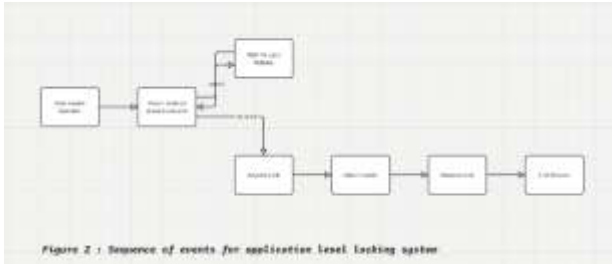
split-brain, where multiple nodes think they are leaders, the developer needs to program a robust detection mechanism for leader, recovery, and consensus.

#### 4. Implementation Overview.

One such unique approach to the election of a leader for an application is by applying distributed lockers or mutex. The primary goal of a leader election algorithm is to elect a single process as the coordinator among distributed nodes, even in the presence of failures (5). In the approach, the nodes present in the system fight to acquire a lock on a shared reference and the first one to prevail in acquiring the lock is elected through the approach. The leader can do the tasks of a leader by one at a time. When the leader node fails or quits the mutex acquired, the rest of the nodes present in the system can battle to acquire the lock on the shared reference. The nodes present in the system can, therefore, always be checking the mutex status so as to retry and tie to acquire the lock on the reference.

To ensure the robustness of the leader election mechanism, node failures should be handled properly. To detect node failures, the leader can send heartbeat signals to all the nodes, which is a signal that the leader is still active and active. In case a node

doesn't receive heartbeat signals within the time frame then it elects a new leader. When trying to acquire a lock in an entity, there should be a timeout value, so that in case something goes wrong, it will not wait indefinitely and after a certain number of retries it can give up. These are some of the precautions that can be taken to avoid being subject to the above scenario. **Figure 2** shows the sequence of events that will occur during the lifetime of leader selection.



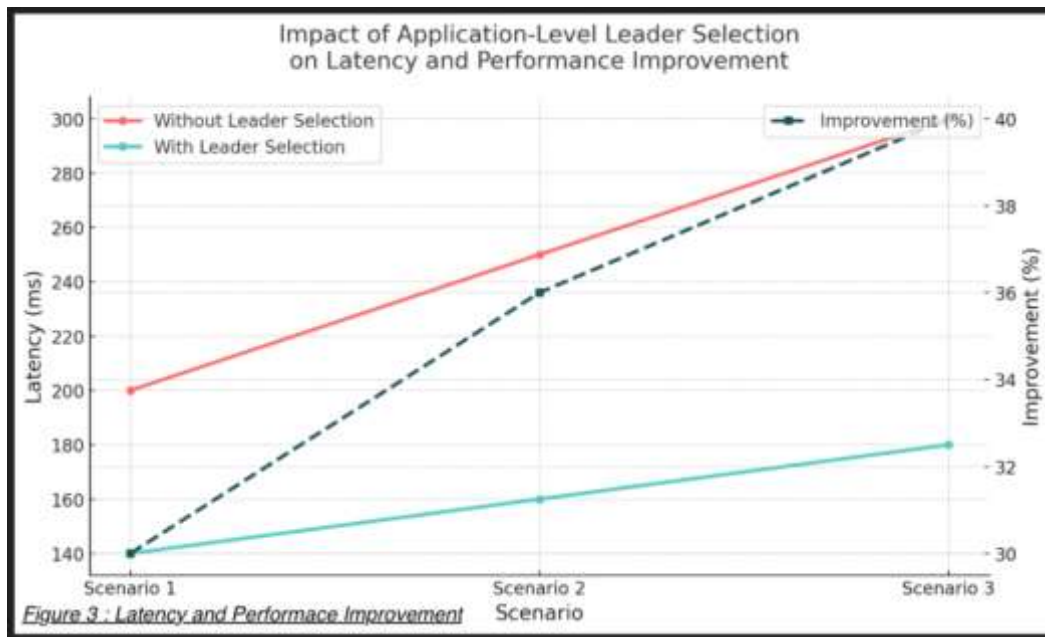
While selecting an application-level leader, it is of utmost importance to bear in mind issues of scalability and performance. As the number of nodes grows, the cost of managing lock acquisition attempts, and failures scales on the upside. This can be mitigated by leveraging techniques such as exponential backoff for retries, as well as sharding the leadership responsibilities across various entities and determining the optimal frequency through which to send heartbeats. Doing this allows you to effectively split the difference inherent in ensuring that the leader that's elected is up, as well as protecting the performance of the system. Leader election is a powerful technique for simplifying distributed architectures by reducing coordination

overhead and ensuring consistent decision making (1).

Careful design of application-level leader election is necessary to avoid pitfalls like a split-brain scenario where multiple nodes think they are the leader. To achieve this, make sure the lock acquisition and release procedures are atomic and the nodes have a similar understanding of lock state. Also, logging leader transitions and monitoring the leader's health can give insights and show what's happening in your distributed system, helping to debug fast enough. By following this process, you can leverage effective and reliable election of leaders without needing coordination services.

### 5. Performance Evaluations.

There are some advantages of application-level leader election, especially in scenarios where the ease of use and the direct relationship with the application are essential. By integrating leader election logic into the application itself, the process of selecting the leader and managing failover can be done with low latency. This approach also eliminates the need for coordination services and simplifies the architecture, which in turn results in lower network overhead and less synchronization. Most importantly, there is complete control over the tuning of election parameters to the application, allowing all of these parameters to be optimized for low latency application-level leader selection and failover. **Figure 3** showcases the benefit of application level leader selection in terms of latency and performance.



Compared to an external coordination service like ZooKeeper or etcd, application-level leader election

would provide lower latency and less operational overhead, as there are no interactions with an

external system. However, the application must now deal with leader election logic and failover. In contrast, a consensus algorithm like Raft or Paxos provides strong consistency, but it may have higher latency and complexity due to multiple rounds of communication. Application-level leader election falls between these two in terms of complexity versus latency and might be suitable when an application can handle leader election effectively without degrading system performance.

Application level leader election is very much beneficial in the cases where application needs some specific requirement for leader selection when setting up the leader election with the help of coordination services is not an easy task. If we consider microservices, a design where a large application has been split into a number of services, having application level leader election means that every service has its own leader election, and it can be modified according to the requirements of the services. Besides, in cases where high latency and low throughput are required, it is in the benefit of the application using application level leader election are very effective due to the reason that they can reduce the overhead produced by coordination services.

## 6. Conclusion

Leader election plays a pivotal role in distributed systems, impacting various applications by ensuring coordination, consistency, and fault tolerance (2). Application layer leader election provides a lightweight and flexible means of managing leadership of distributed systems, especially if external coordination service is not suitable or available. By having the leader election logic built within the application, the system can reduce latency and simplify its architecture. However, it means that the application is now responsible for managing leader election, failure detection, and recovery mechanisms, which can possibly lead to performance bottlenecks if not implemented properly. Leader election algorithms form the backbone of many distributed systems, providing coordination, fault tolerance, and consistency through structured leadership roles (6).

To have stronger and more scalable application-level leader election, future work could concentrate on incorporating adaptive approaches that regulate leader election parameters based on system performance. For example, designing performance-oblivious leader election algorithms such as SEER which is able to determine the replica with the lowest average cluster response time, could lead to an effective selection of the leader in different workloads and network conditions.

In addition, finding a hybrid approach to combine application level leader election and lightweight some coordination services can provide autonomy and reliability. For example, you can use a lease based locking which is a service as ZooKeeper, DynamoDb and therefore, has some failover capabilities and fault tolerance and still be simple like application level leader election.

Similarly, there is a need to address how network partitions and node crashes will be handled. This can be achieved through mechanisms such as quorum-based voting or using timeouts and retries to ensure that leader election remains consistent even in adverse conditions. Research into partially asynchronous leader election algorithms such as PALE would be of benefit as they operate effectively in dynamic systems with weak assumptions about message delays and clock drift.

In conclusion, the program-level leader election proposes a clear and reliable leader management methodology for distributed systems, but it is possible for it to be improved considerably. By incorporating adaptive algorithms, hybrid coordination strategies, and strong failure mechanisms, the performance and reliability of the program-level leader election can be significantly increased, making it a better option for large-scale and complex distributed systems.

## Author Statements:

- **Ethical approval:** The conducted research is not related to either human or animal use.
- **Conflict of interest:** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper
- **Acknowledgement:** The authors declare that they have nobody or no-company to acknowledge.
- **Author contributions:** The authors declare that they have equal right on this paper.
- **Funding information:** The authors declare that there is no funding to be acknowledged.
- **Data availability statement:** The data that support the findings of this study are available on request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions.

## References

- [1] Amazon Builders' Library – Leader Election in Distributed Systems - <https://aws.amazon.com/builders-library/leader-election-in-distributed-systems/>

- [2] UNM Research: Scalable Leader Election. - <https://www.cs.unm.edu/~saia/papers/leader.pdf>
- [3] PingCAP: Innovative Leader Election in Distributed Systems with S3 Writes.- <https://www.pingcap.com/article/innovative-leader-election-in-distributed-systems-with-s3-writes/>
- [4] Medium – Leader Election in Distributed Systems - <https://medium.com/@sutanu3011/leader-election-in-distributed-system-8f3d746853c0>
- [5] GeeksforGeeks – Leader Election in System Design - <https://www.geeksforgeeks.org/leader-election-in-system-design/>
- [6] ResearchGate – A Survey and Taxonomy of Leader Election Algorithms - [https://www.researchgate.net/publication/342927985\\_A\\_Survey\\_and\\_Taxonomy\\_of\\_Leader\\_Election\\_Algorithms\\_in\\_Distributed\\_Systems](https://www.researchgate.net/publication/342927985_A_Survey_and_Taxonomy_of_Leader_Election_Algorithms_in_Distributed_Systems)