

## Exploration of Distributed Relational Databases in the Internet Era

Xie Zhixiao  
Nanning College, Nanning Guangxi

---

**Abstract:** In the internet era, the scale and complexity of data have increased dramatically, making traditional centralized databases inadequate for handling the demands of modern applications. Distributed relational databases offer a solution with their ability to manage large volumes of data across multiple locations, ensuring scalability, availability, and performance. This paper examines the architecture, benefits, and challenges of distributed relational databases, including data consistency, integrity, and the support for contemporary data-driven applications.

**Keywords:** distributed relational databases; data scalability; data consistency; internet era; data integrity

---

### I. The Emergence of Distributed Relational Databases

#### A. Background and Necessity in the Internet Era

The internet era has been characterized by an explosion of data, driven by the proliferation of digital devices, the internet of things (IoT), and the increasing digitalization of services and industries. This data growth has outpaced the capabilities of traditional centralized databases, which struggle to scale horizontally to accommodate larger datasets and higher query loads. The centralized model, with all data stored in a single location, also faces challenges in terms of availability and disaster recovery. In this context, the necessity for distributed relational databases has become evident.

Distributed databases distribute data across multiple nodes, which can be located in different physical locations, thereby providing better scalability and resilience. They are designed to handle large volumes of structured data while maintaining the relational integrity and the ability to execute complex queries, which are essential for transactional systems and data analytics.

#### B. Comparison with Traditional Centralized Databases

In contrast to centralized databases, distributed relational databases offer several key advantages. Centralized databases rely on a single server or a limited number of servers to store and process data. While

this model is simple and easy to manage, it becomes a bottleneck as data grows and as the demand for data access increases. On the other hand, distributed databases distribute the load across multiple nodes, which can work in parallel to process queries and transactions.

One of the main challenges with centralized databases is their susceptibility to single points of failure. If the central server goes down, the entire database becomes inaccessible. Distributed databases mitigate this risk by replicating data across different nodes, ensuring that the system remains available even if one or more nodes fail.

Another significant difference is the ability to scale out. Centralized databases can scale vertically by upgrading hardware, but this approach has limitations and can be costly. Distributed databases, in contrast, can easily scale out by adding more nodes to the system, which is a more flexible and cost-effective solution to increasing data and query loads.

#### C. Core Characteristics and Advantages

Distributed relational databases are defined by several core characteristics that give them an edge over traditional centralized systems:

1. **Scalability:** Distributed databases can scale out horizontally by adding more nodes to the system, allowing them to handle increasing amounts of data and

user requests.

2. **Availability:** Through data replication, distributed databases ensure high availability. Even if some nodes are unavailable, the system can continue to function, providing access to the data.

3. **Fault Tolerance:** The distributed nature of these databases means that they can tolerate node failures without losing data or stopping the service.

4. **Performance:** By distributing the workload, distributed databases can improve performance, as read and write operations can be parallelized across multiple nodes.

5. **Consistency:** Distributed databases employ various consistency models to ensure that data remains accurate and up-to-date across all nodes, despite the challenges posed by replication latency.

6. **Flexibility:** They offer flexibility in terms of data placement and workload management, allowing administrators to optimize the system based on the specific access patterns and data characteristics.

7. **Geo-Distribution:** For global applications, distributed databases can be deployed across multiple geographical locations, reducing latency for local users and providing disaster recovery capabilities.

8. **Cost-Effectiveness:** The ability to use commodity hardware for the database nodes makes distributed databases a cost-effective solution compared to the high-cost, vertically scalable systems.

In conclusion, the emergence of distributed relational databases is a response to the challenges posed by the internet era's data demands. Their design addresses the limitations of centralized databases and provides a robust foundation for modern applications that require high scalability, availability, and performance. As the volume and velocity of data continue to grow, the advantages of distributed relational databases will become increasingly vital for organizations relying on data to drive their operations and strategies.

## II. Design and Operation of Distributed Systems

### A. Architectural Considerations for Distribution

The design of distributed relational databases

involves several key architectural considerations. The primary goal is to create a system that is not only scalable and fault-tolerant but also capable of maintaining the ACID (Atomicity, Consistency, Isolation, Durability) properties of transactions.

1. **Decentralization:** Distributed databases move away from the centralized control model, distributing both data and control logic across the network. This decentralization requires sophisticated mechanisms for coordination and data management.

2. **Network Topology:** The arrangement of nodes in a distributed system can significantly affect performance and reliability. Common topologies include star, ring, and fully meshed networks, each with its advantages and trade-offs.

3. **Data Locality:** Designing for data locality, where data is stored close to the processing power that manipulates it, can reduce latency and network traffic.

4. **Load Balancing:** Effective distribution of workloads across nodes is crucial for optimizing resource utilization and response times.

5. **Scalability:** The system must be designed to scale out by adding more nodes as demand increases, without significant reconfiguration.

6. **Interoperability:** Ensuring that the distributed system can work seamlessly with existing applications and services is essential for adoption and integration.

### B. Data Partitioning and Replication Strategies

Data partitioning and replication are fundamental strategies in distributed databases that contribute to their performance and reliability.

1. **Horizontal Partitioning (Sharding):** Involves splitting tables into smaller, manageable pieces based on a shard key, distributing the load and improving query performance.

2. **Vertical Partitioning:** Divides tables by functionality, storing frequently accessed columns on faster storage media to optimize read operations.

3. **Replication:** Duplicates data across multiple nodes to ensure high availability and durability. Replication

can be synchronous or asynchronous, depending on the required balance between consistency and performance.

4. **Consistency Models:** Different replication strategies lead to various consistency models, such as strong consistency, eventual consistency, and causal consistency, each with its use cases and implications for application design.

5. **Conflict Resolution:** In environments with asynchronous replication, mechanisms for conflict detection and resolution are necessary to handle data discrepancies across replicas.

### C. Ensuring Consistency and Synchronization

Consistency and synchronization are central challenges in distributed systems.

1. **Distributed Transactions:** Ensuring that a set of operations across multiple nodes is treated as a single atomic unit is complex. The two-phase commit protocol is a traditional method for achieving this but can be a bottleneck.

2. **Consensus Algorithms:** Algorithms like Raft or Paxos are used to reach consensus across distributed nodes, which is essential for maintaining consistency in the absence of a central authority.

3. **Timestamp Ordering:** Some distributed databases use a logical timestamp ordering to maintain consistency without the need for locks, allowing for higher concurrency.

4. **Vector Clocks:** A data structure used to track causality and concurrency in distributed systems, helping to resolve conflicts and maintain consistency.

5. **Synchronization Protocols:** Protocols for data synchronization ensure that replicas are kept up-to-date. These can range from simple periodic synchronization to more complex event-driven or change-data-capture (CDC) based approaches.

6. **Monitoring and Alerts:** Implementing robust monitoring to detect and respond to synchronization issues and performance bottlenecks is critical for maintaining system health.

In summary, the design and operation of distributed

systems require careful consideration of architectural decisions that impact performance, consistency, and reliability. Data partitioning and replication strategies are employed to enhance these aspects, while maintaining strict consistency and synchronization across the distributed environment is an ongoing challenge that demands sophisticated solutions. As distributed relational databases continue to evolve, these design considerations will remain central to their success in managing the vast and dynamic datasets of the internet era.

## III. Data Integrity and Transaction Management

Data integrity and transaction management are critical components of any robust database system. Ensuring that data remains consistent, accurate, and reliable is paramount, especially in distributed environments where data is spread across multiple nodes. This section will delve into the importance of preserving ACID properties in distributed environments, the challenges that arise in managing distributed transactions, and the various solutions and protocols designed to maintain data integrity.

### A. Preserving ACID Properties in Distributed Environments

ACID (Atomicity, Consistency, Isolation, Durability) properties are the cornerstone of transaction processing. These properties ensure that database transactions are processed reliably and consistently. In a distributed environment, preserving ACID properties becomes more challenging due to the inherent complexities of managing data across multiple nodes. Let's examine each property and its significance in distributed environments:

1. **Atomicity:** Atomicity ensures that a transaction is treated as a single, indivisible unit of work. In a distributed environment, atomicity can be challenging to maintain because a transaction may involve multiple nodes. If one part of the transaction fails, the entire transaction must be rolled back to maintain atomicity. Techniques such as two-phase commit (2PC) and three-phase commit (3PC) protocols are used to coordinate the commit or abort of distributed transactions.

2. **Consistency:** Consistency ensures that a transaction brings the database from one valid state to another. In a distributed environment, maintaining consistency is crucial because data is replicated across multiple nodes. Any changes made to the data must be propagated to all replicas to ensure consistency. Techniques like distributed consensus algorithms (e.g., Paxos, Raft) help achieve consistency by ensuring that all nodes agree on the state of the data.

3. **Isolation:** Isolation prevents interference between concurrent transactions. In a distributed environment, isolation becomes more complex due to the potential for multiple transactions to access the same data simultaneously across different nodes. Isolation levels, such as serializable and snapshot isolation, are used to control the interaction between concurrent transactions and maintain data integrity.

4. **Durability:** Durability ensures that once a transaction is committed, its effects are permanent and will survive any subsequent system failures. In a distributed environment, durability can be achieved through techniques like write-ahead logging (WAL) and replication. WAL ensures that changes are logged before they are applied to the database, while replication ensures that data is copied to multiple nodes for fault tolerance.

## B. Challenges in Managing Distributed Transactions

Managing distributed transactions presents several challenges that must be addressed to ensure data integrity and system performance. Some of these challenges include:

1. **Network Latency and Partitioning:** In a distributed environment, network latency and partitioning can occur, leading to communication delays and potential data inconsistencies. Ensuring that transactions are executed correctly despite these issues requires sophisticated algorithms and protocols.

2. **Scalability:** As the number of nodes in a distributed system increases, scalability becomes a concern. The system must be able to handle a growing number of transactions and data without sacrificing

performance or data integrity.

3. **Concurrency Control:** Managing concurrency in a distributed environment is more complex than in a single-node system. Ensuring that transactions execute concurrently without causing conflicts or inconsistencies requires advanced concurrency control mechanisms.

4. **Data Replication and Synchronization:** Data replication is essential for fault tolerance and performance, but it also introduces challenges in maintaining consistency across replicas. Synchronizing data across multiple nodes can be difficult, especially in the presence of network partitions and node failures.

5. **Transaction Coordination:** Coordinating distributed transactions requires careful management to ensure that all participating nodes agree on the outcome of the transaction. Protocols like 2PC and 3PC can be complex to implement and may introduce additional overhead.

## C. Solutions and Protocols for Data Integrity

To address the challenges of managing distributed transactions and preserving data integrity, several solutions and protocols have been developed. Some of these include:

1. **Two-Phase Commit (2PC):** The 2PC protocol is a widely used method for achieving atomicity in distributed transactions. It involves a coordinator node that manages the commit process by collecting votes from participating nodes and deciding whether to commit or abort the transaction.

2. **Three-Phase Commit (3PC):** The 3PC protocol is an extension of 2PC that aims to reduce the time a coordinator is locked. It introduces an additional phase to improve the chances of reaching a consensus and committing the transaction.

3. **Distributed Consensus Algorithms:** Algorithms like Paxos and Raft are used to achieve consensus among nodes in a distributed system. These algorithms ensure that all nodes agree on the state of the data, which is crucial for maintaining consistency.

4. **Replication and Synchronization Techniques:**



Techniques like multi-master replication, where multiple nodes can accept write operations, and eventual consistency, where data consistency is achieved over time, are used to manage data replication and synchronization in distributed systems.

5. Concurrency Control Mechanisms: Various concurrency control mechanisms, such as optimistic and pessimistic locking, are used to manage concurrent access to data in distributed environments. These mechanisms

#### IV. Applications and Future Prospects

Distributed databases have become an integral part of modern data management solutions, enabling a wide range of applications and driving innovation in various industries. This section explores the role of distributed databases in enabling big data and real-time analytics, their significance in cloud computing and microservices architectures, and the emerging trends that are shaping the future of distributed databases.

##### A. Enabling Big Data and Real-time Analytics

The exponential growth of data in recent years has led to the concept of big data, which refers to the vast volumes of data generated at high velocity and in diverse formats. Distributed databases are uniquely positioned to handle the challenges posed by big data due to their scalability, performance, and ability to process large datasets in parallel.

1. Scalability: Distributed databases can scale horizontally by adding more nodes to the system, which allows them to handle increasing data volumes without sacrificing performance. This scalability is crucial for big data applications that require processing and storing petabytes of data.

2. Real-time Analytics: The distributed nature of these databases allows for real-time processing of data, enabling organizations to gain immediate insights and make data-driven decisions. Real-time analytics is particularly valuable in fields like finance, healthcare, and IoT, where timely information can lead to significant advantages.

3. Data Variety: Distributed databases are designed

to handle diverse data types, including structured, semi-structured, and unstructured data. This flexibility is essential for big data applications that need to integrate and analyze data from various sources.

##### B. Role in Cloud Computing and Microservices Architectures

Cloud computing and microservices architectures have revolutionized the way applications are developed and deployed. Distributed databases play a pivotal role in these modern computing paradigms.

1. Cloud Computing: Cloud-based distributed databases offer elasticity, allowing resources to be scaled up or down based on demand. They also provide high availability and fault tolerance, which are critical for cloud services. Cloud providers like Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP) offer managed distributed database services that simplify deployment and management.

2. Microservices Architectures: In microservices architectures, applications are composed of small, independent services that can be developed, deployed, and scaled separately. Distributed databases are well-suited for microservices because they can provide each service with its own data store while maintaining data consistency across the system.

3. Service Discovery and Orchestration: Distributed databases can be integrated with service discovery and orchestration tools like Kubernetes to manage the lifecycle of microservices, ensuring that services can find and communicate with each other efficiently.

##### C. Emerging Trends and the Road Ahead for Distributed Databases

The field of distributed databases is rapidly evolving, with several trends shaping its future:

1. NewSQL and Hybrid Databases: NewSQL databases aim to combine the scalability of NoSQL systems with the ACID guarantees of traditional SQL databases. Hybrid databases that offer both transactional and analytical capabilities are also gaining traction, providing a single platform for various data workloads.

2. **Serverless Databases:** Serverless architectures abstract away the infrastructure management, allowing developers to focus on application logic. Serverless databases are emerging as a trend, where the database service automatically scales in response to the application's needs without the need for manual provisioning.

3. **Edge Computing:** With the rise of IoT and the need for real-time processing at the edge of the network, distributed databases are being deployed closer to the data sources. Edge databases can reduce latency and handle data processing tasks locally.

4. **Machine Learning Integration:** Distributed databases are increasingly being integrated with machine learning frameworks to enable in-database machine learning. This integration allows for more efficient processing of large datasets and simplifies the deployment of machine learning models.

5. **Blockchain and Distributed Ledgers:** Blockchain technology, which relies on distributed databases, is being explored for applications beyond cryptocurrencies. It offers a decentralized and secure way to manage transactions and could have significant implications for data integrity and trust in various industries.

6. **Data Privacy and Security:** As data privacy regulations become more stringent, distributed databases must evolve to include advanced encryption, access controls, and compliance features to protect sensitive data.

## V. Conclusion

### A. Summary of Key Points and Findings

This paper has explored the emergence and operational dynamics of distributed relational databases in the internet era. We began by examining the necessity of distributed systems due to the exponential growth of data and the limitations of traditional centralized databases. The comparison revealed that distributed databases offer superior scalability, availability, and performance, which are critical for modern applications.

Key characteristics such as decentralization, data locality, and load balancing were identified as the

cornerstones of distributed database architecture. The paper also discussed the importance of data partitioning and replication strategies, which are essential for managing large datasets and ensuring high availability and fault tolerance.

Furthermore, we delved into the complexities of ensuring consistency and synchronization in a distributed environment. The discussion highlighted the role of distributed transactions, consensus algorithms, and synchronization protocols in maintaining data integrity and system reliability.

### B. Significance of Distributed Relational Databases in Modern Infrastructure

Distributed relational databases have become a linchpin in modern data infrastructure. They provide the backbone for applications ranging from e-commerce platforms and social networks to big data analytics and IoT systems. The ability to handle large volumes of data with low latency and high throughput is crucial for real-time processing and decision-making.

Moreover, the distributed nature of these databases aligns with the geographic distribution of users and services in the globalized digital economy. By reducing latency through geo-distribution and replicating data across regions, distributed databases enhance the user experience and support disaster recovery strategies.

The significance of distributed relational databases is further amplified by their role in facilitating cloud computing and microservices architectures. They offer the flexibility needed to support dynamic scaling and the diverse data requirements of cloud-based applications.

### C. The Imperative for Ongoing Research and Development

Despite the advancements in distributed relational databases, there is a clear imperative for ongoing research and development. Emerging technologies such as artificial intelligence, machine learning, and advanced analytics are driving new demands for data management capabilities.

Challenges related to data security, privacy, and compliance in a distributed context also require innovative

solutions. As data becomes more decentralized, ensuring data protection and meeting regulatory requirements become increasingly complex.

The future of distributed relational databases will likely see more integration with emerging technologies, the development of more sophisticated consistency models, and enhanced support for global, real-time applications. It is crucial for the research community and industry to collaborate in driving these innovations forward.

## References

- [1] Yang Zhenkun, Yang Chuanhui, Han Fusheng, et al. OceanBase Distributed relational database architecture and technology [J]. Computer Research and Development, 2024,61 (3): 540-554.
- [2] Zheng Xiaofeng, Zhou Ke. Implementation of a distributed relational database interconnection based on OptiConnect [J]. Journal of Jiangsu Normal University of Technology, 2003,9 (4): 55-60.
- [3] Zheng Xiaofeng, Ma Zhenghua. Application of a distributed relational database architecture under the AS / 400 cluster [J]. Journal of Jiangsu Institute of Technology, 2003,15 (1): 57-60.
- [4] Li Yigang. Keyword query [D] on a distributed relational database. Heilongjiang Province: Harbin Institute of Technology, 2009.
- [5] Yu Jun. TiDB new generation distributed relational database financial industry innovation practice [C]. / Proceedings of the 2018 DAMS China Data Asset Management Summit. 2018:1-26.
- [6] Jiang Mingjun. The Design and Implementation of the Distributed Relational Database Transaction Manager [D]. Jiangsu: Southeast University, 2019.
- [7] Dai Zhengbing. Oracle Application and research of distributed relational database in FXCMIS system [D]. Hunan: Hunan University, 1998.
- [8] Yang Zhenkun. Commercial reliability of the distributed relational database OceanBase [J]. Electronic Finance, 2016,0 (2): 67-69.
- [9] Zhou Qian Jun. Application of the distributed relational database ORACLE in LDMIS [D]. Hunan: Hunan University, 1996.
- [10] Lu Hanrong. Dynamic data duplication for a distributed relational database [J]. Computer Engineering, 1994, (4): 11-13,27.
- [11] Sun Bin, Feng Naiqin. Distributed database incremental update method based on the relational model [J]. Computer Simulation, 2024,41 (5): 518-521,531.
- [12] Yang Zhenkun. High reliability of the distributed relational database OceanBase [J]. Electronic Finance, 2016 (2): 67-69.
- [13] State Administration of Quality Supervision, Inspection and Quarantine of the People's Republic of China. Distributed Relational Database service interface specification: GB / T 32633-2016 [S]. 2016.
- [14] Han Ye. Research on the distributed storage method of non-relational database based on big data technology [J]. Information and Computers, 2024,36 (4): 166-168.

©2024. This article is copyrighted by the author and Hong Kong Science and Technology Publishing Group. This work is licensed under a Creative Commons Attribution 4.0 International License.

<http://creativecommons.org/licenses/by/4.0/>



**Open Access**