Lecture 3

# **Basic Principles: CAP theorem, Consistency**

Yuliia Prokop prokoyul@fel.cvut.cz

6. 10. 2025

Based on the presentation of Martin Svoboda (martin.svoboda@matfyz.cuni.cz)

THE STATE OF THE S



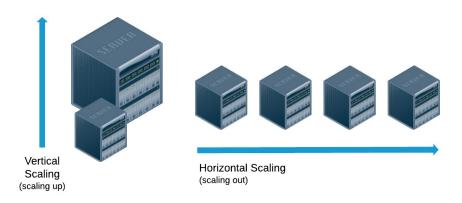
Czech Technical University in Prague, Faculty of Electrical Engineering

### **Lecture Outline**

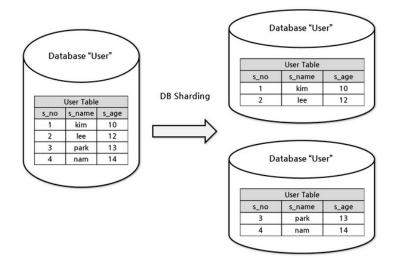
#### Different aspects of data distribution

- Scaling
  - Vertical vs. horizontal
- **Distribution** models
  - Sharding
  - Replication: master-slave vs. peer-to-peer architectures
- CAP properties
  - Consistency, availability and partition tolerance
  - ACID vs. BASE guarantees
- Consistency
  - Read and write quorums

# Lecture 2 overview: Horizontal vs. Vertical Scaling



# **Lecture 2 overview: Sharding**



# Replication

# Replication in distributed systems

**Replication** refers to maintaining identical copies of data across multiple servers.

#### The primary motivations include:

Fault tolerance:

System resilience against individual node failures

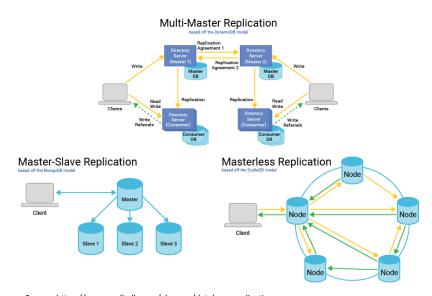
High availability:

Continued operation during partial system failures

Performance optimization:

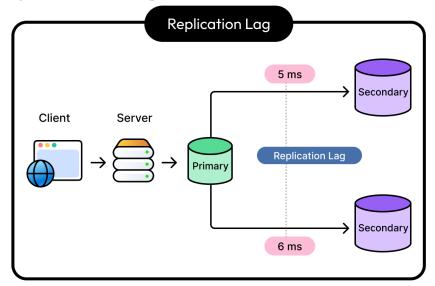
Reduced latency through geographic proximity

### Replication in distributed systems



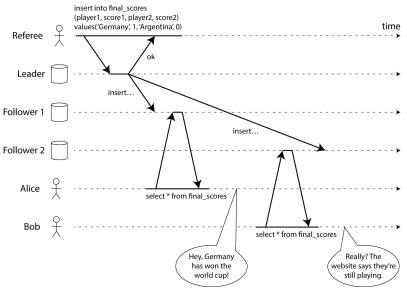
Source: https://www.scylladb.com/glossary/database-replication

# **Replication Lag**



Source: hhttps://blog.bytebytego.com/p/a-guide-to-database-replication-key

# **Replication Lag**



Source: https://tarunjain07.medium.com/cap-theorem-notes-68b04523cbce

### **Consensus Protocols**

Consensus protocols ensure data consistency among replicas. They are essential for making unified decisions about the system's state, especially during failures and network partitions.

#### **Main Protocols:**

- Paxos
- Raft
- Viewstamped Replication
- Zookeeper's Zab Protocol.

#### **Role of Consensus Protocols:**

- Ensuring data consistency across replicas.
- Leader or master election within the system.
- Coordinating data updates and managing node failures.

### **CAP Theorem**

### **CAP Theorem**

#### Assumptions

- Distributed system with sharding and replication
- Read and write operations on a single aggregate only

#### **CAP** properties

- Properties of a distributed system
- Consistency, <u>A</u>vailability, and <u>P</u>artition tolerance CAP theorem

In the presence of a network **partition**, a distributed system can choose either **consistency** or **availability**, but not both.

But, what these properties actually mean?

#### C: Consistency

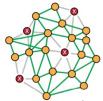
At any given time, all nodes in the network have exactly the same (most recent) value.



= Value: X @ 2018-05-03T08:52:40

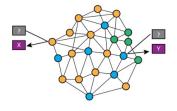
#### P: Partition tolerance

The network continues to operate, even if an arbitrary number of nodes are failing.



#### A: Availability

Every request to the network receives a response, though without any guarantee that returned data is the most recent.



- O = Value: X @ 2018-05-03T08:52:40
- = Value: Z @ 2018-05-03T08:32:58
- = Value: Y @ 2018-05-03T07:12:12

Source: https://tarunjain07.medium.com/cap-theorem-notes-68b04523cbce

# **CAP Properties**

Property	Formal Definition	Practical Meaning	
Consistency	Linearizability: Operations appear to execute atomically	All reads return the most recent write	
Availability	Every request receives a response (success or failure)	The system always responds, never times out	
Partition Tolerance	System continues despite message loss between nodes	Works even when network splits occur	

- Hardware failures are inevitable
- Network congestion causes effective partitions
- Slow networks trigger timeouts
- Geographic distribution increases partition probability

# **CAP Properties: Consistency Details**

- Intuition: Each read/write on a key is atomic.
- Formal (linearizability): There is a single global order of operations; each operation takes effect at an instantaneous point between its call and completion as if all ran sequentially on one node.
- **Consequence:** After a successful write, any later read of the same key returns the updated value (read-after-write).
- **Replication requirement:** If any replica can serve reads, a write must be replicated to a sufficient set (e.g., a quorum) before acknowledgment to preserve strong consistency.
- Weaker consistency models also exist.

### **CAP Properties**

### **Availability**

- If a node is working, it must respond to user requests
  - A bit more formally...

Every read or write request successfully <u>received</u> by a non-failing node in the system must result in a response (success or failure), not be silently dropped.

I.e., their execution must not be rejected

### **Partition tolerance**

 The system continues to operate even when two or more sets of nodes get isolated

A bit more formally...

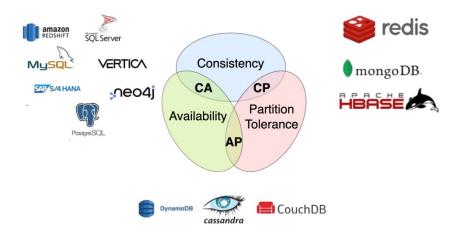
The network is allowed to lose arbitrarily many messages sent from one node to another

I.e. a connection failure must not shut the whole system down

### **CAP Theorem Proof**

- Proof by contradiction
  - Assume all three properties can be satisfied simultaneously
  - Consider a network partition scenario
- Partition scenario setup
  - Network splits into two disjoint sets of nodes: G<sub>1</sub> and G<sub>2</sub>
  - No communication possible between G₁ and G₂
- Write operation on G<sub>1</sub>
  - Client writes to G<sub>1</sub>, must be consistent across all replicas
  - G<sub>2</sub> cannot receive this update due to partition
- Read operation on G<sub>2</sub>
  - If system is available, G<sub>2</sub> must respond to read requests
  - If system is consistent: G<sub>2</sub> must return the updated value
- ✓ Contradiction: G₂ cannot have updated value (violates C) but must respond (requires A)

### $C \wedge A \wedge P$ is impossible in distributed systems



# **Consistency Spectrum**

### Strong consistency models

- Linearizability (strongest for a single operation/key)
- Transactional models (Serializability / Snapshot Isolation)
- Sequential consistency
- Causal consistency

### Weak consistency models

- Session consistency
- Monotonic read/write consistency
- Eventual consistency (weakest)

### Consistency vs. Performance trade-off

- Stronger consistency → Higher latency
- Weaker consistency → Better performance

### Application requirements determine choice

- Banking: Strong consistency required
- Social media: Eventual consistency acceptable
- Collaborative editing: Causal consistency needed

If at most two properties can be guaranteed...

- CA = consistency + availability
  - Traditional ACID properties are easy to achieve
  - Examples: RDBMS
  - Any single-node system, but even clusters (at least in theory)
    - However, should the network partition happen, all the nodes must be forced to stop accepting user requests

CA: Consistency + Availability – only possible if no network partitions occur

(e.g., traditional RDBMS under normal conditions)

If at most two properties can be guaranteed...

- CP = consistency + partition tolerance
  - Other examples: distributed locking
- AP = availability + partition tolerance
  - New concept of BASE properties
  - Examples: Apache Cassandra, Apache CouchDB.
  - Other examples: web caching, DNS

In real-world environments, network partitions can and do occur. Distributed systems therefore **should be designed to tolerate partitions (P)** and then choose between C and A during a partition. Systems that sacrifice P effectively stop responding when a partition occurs.

#### Design for partitions in clusters

- Why?
  - Because it is difficult to detect network failures
- Does this mean that only purely CP and AP systems are possible?
  - No...

#### **The real meaning** of the CAP theorem:

- The real world does not need to be just black and white
- Partition tolerance is a must, but we can trade off consistency versus availability
  - A relaxed consistency can bring a lot of availability.
  - Such trade-offs are not only possible, but often work very well in practice

# **ACID Properties**

#### Traditional **ACID** properties

- Atomicity
  - Partial execution of transactions is not allowed (all or nothing)
- Consistency
  - Transactions bring the database from one consistent (valid) state to another
- Isolation
  - Transactions executed in parallel do not see uncommitted effects of each other
- Durability
  - Effects of committed transactions must remain durable

# **BASE Properties**

#### New concept of **BASE** properties

- <u>B</u>asically <u>A</u>vailable
  - The system works basically all the time
  - Partial failures can occur, but there are no total system failures
- Soft State
  - The system is in flux (unstable), non-deterministic state
  - Changes occur all the time
- Eventual Consistency
  - Sooner or later the system will be in some consistent state

BASE is just a vague term, no formal definition was provided

 Proposed to illustrate design philosophies at the opposite ends of the consistency-availability spectrum

### **ACID** and **BASE**

#### **ACID**

- Choose <u>consistency over availability</u>
- Pessimistic approach
- Implemented by traditional relational databases

#### **BASE**

- Choose <u>availability over consistency</u>
- Optimistic approach
- Common in NoSQL databases
- Allows levels of scalability that cannot be acquired with ACID

Historical move:

strong consistency → eventual consistency

Current trend in NoSQL:

eventual only → tunable/stronger consistency options

# **Consistency**

# **Consistency**

Consistency in general...

- Consistency is the lack of contradiction in the database
- However, it has many facets...
  - For example, we only assume atomic operations that constantly manipulate a single aggregate.
     But set operations could also be considered, etc.

Strong consistency is achievable in clusters with appropriate replication/consensus (e.g., quorum/majority, consensus protocols), but eventual consistency might often be sufficient.

- A one-minute-old article on a news portal does not matter
- Even when an already unavailable hotel room is booked once again, the situation can still be figured out in the real world

• ...

### **Consistency vs. Latency Trade-offs**

### Strong consistency costs

- Synchronous replication to a quorum/majority of nodes
- Latency ≈ latency to the slowest node in the quorum
- Example: 3 nodes, majority = 2, 100 ms max → ~100 ms latency

#### Weak consistency benefits

- Asynchronous replication
- Latency = latency to a single node
- Example: 3 nodes, 10ms local → 10ms total latency

#### Real-world measurements

- MongoDB: 5ms local read, 50ms strongly consistent read
- Cassandra: 2ms eventual read, 20ms quorum read

### Tunable consistency (modern approach)

- Applications can choose per-operation
- Critical operations: strong consistency
- Non-critical operations: eventual consistency

### Consistency

### Write consistency (update consistency)

- Problem: write-write conflict
  - Two or more write requests on the same aggregate are initiated concurrently
- Context: peer-to-peer architecture only
- Issue: lost update
- Solution:
  - Pessimistic strategies
    - Preventing conflicts from occurring
    - Write locks, ...
  - Optimistic strategies
    - Conflicts may occur, but are detected and resolved later on
    - Version stamps, vector clocks, ...

# **Consistency**

#### Read consistency (replication consistency)

- Problem: read-write conflict
  - Write and read requests on the same aggregate are initiated concurrently
- Context: both master-slave and peer-to-peer architectures
- Issue: inconsistent read
- When not treated, inconsistency window will exist
  - Propagation of changes to all the replicas takes some time
  - Until this process is finished, inconsistent reads may happen
  - Even the initiator of the write request may read wrong data!
    - Session consistency / read-your-writes / sticky session

### **Strong Consistency**

How many nodes need to be involved to get strong consistency?

**General rule:** R + W > N (read and write quorums must intersect)

- Write quorum: W > N/2
  - Idea: a majority write ensures only one write can succeed at a time
    - W = number of nodes successfully acknowledged the write
    - N = number of nodes involved in replication (replication factor)
- Read quorum: choose R such that R + W > N (e.g., R > N W)

Idea: intersecting quorums ensure reads see the latest committed write R =number of nodes participating in the read

If the retrieved replicas return different versions, resolve to the **latest committed version** (e.g., via version/timestamp) and then return it.

When a quorum is not attained → the request cannot be handled

# **Strong Consistency**

#### **Examples**

### Examples for replication factor N = 3

- Write quorum W = 3 and read quorum R = 1
  - All the replicas are always updated
  - ullet  $\Rightarrow$  we can read any one of them
- Write quorum W=2 and read quorum R=2
  - Typical configuration, reasonable trade-off

#### Consequence

- Quora can be configured to balance read and write workload
  - The higher the write quorum is required, the lower the read quorum can then be required

### **Bank:**

### **Different Tasks = Different Decisions**

#### Prefer CP semantics

- Account Balance
- Money Transfers
- Loan Approvals
- Transaction Processing
- Credit Limits

#### **Prefer AP** semantics

- Transaction History
- Product Recommendations
- Market News
- Branch Locator
- Customer Chat

### **Online Store: Customer Journey**

### **E-commerce System**

Product
Browsing
AP
Discovery
over accuracy

Cart Mixed Session consistency Check
CP
Prevent
overselling

Payment
Processing
CP
Financial
accuracy

Order
Confirm
CP
Customer
trust

### **University: Academic vs Administrative**

#### **Academic Functions (CP)**

- Student Grades
- Course Registration
- Tuition Payments
- Financial Aid
- Transcripts

### Campus Services (AP)

- Library Search
- Campus Events
- Dining Menus
- Student Organizations
- News & Updates

# University: Critical Example – Course Registration

### **Problem: Popular Course with Limited Seats**

'Machine Learning 101' - 30 seats, 200 students at 8 AM  $\rightarrow$  Need fair, accurate registration

**Solution: CP (Consistency Required):** the system may sacrifice availability to avoid overbooking.

Trade-off: System slower during peak times, but zero overbooking

### **Universal Patterns Across Industries**

Function Type	Bank	E-commerce	University	Pattern
Money/Financial	СР	СР	СР	Usually CP
User Identity	СР	Mixed	СР	Usually CP
Limited Resources	_	СР	СР	Usually CP
Content/Search	AP	AP	AP	Usually AP
History/Logs	AP	AP	AP	Usually AP
Recommendations	AP	AP	AP	Usually AP

Function type predicts CP/AP choice across all industries

### **How to Decide: CP or AP?**

Identify Function Type

Financial? → Usually CP
Content? → Usually AP
Registration? → Usually CP

2 Analyze Error Impact

Money lost? → CP required User frustration? → AP better Legal issue? → CP required

3 User Expectations

Instant response? → AP
Accuracy critical? → CP
Both needed? → Hybrid

4 Design Implementation

CP: Transactions, locks
AP: Caches, replicas
Mixed: Different DBs

### **Lecture Conclusion**

There is a wide range of options influencing...

- Availability when nodes may refuse to handle user requests?
- Consistency what level of consistency is required?
- Latency how long does it take to handle user requests?
- Durability is the committed data written reliably?
- Resilience can the data be recovered in case of failures?

⇒ it's good to know these properties and choose the right trade-off