# Architecture of Distributed Systems

Replication 2019-2020 Consistency

R. H. Mak





Technische Universiteit **Eindhoven** University of Technology

Where innovation starts

# Agenda

- Introduction
  - Partitioning versus replication
  - Quality drivers
  - Replication transparency
  - Architectural concerns
  - Basic architecture
- Consistency models
- Replica management
- Architectural case studies





# Partitioning versus replication

### Partitioning

- Splitting the data into smaller subsets (called partitions) that are distributed over multiple nodes
- Also called sharding

### Replication

- Distributing multiple copies (replicas) of the same data over distinct nodes
- Redundancy

### **Often combined**

• Multiple replicas per partition





# **Quality drivers for partitioning**

Introducing partitioning into the architecture of a data store is mainly motivated by the following quality drivers

- Size of the data set (data base) simply too large for a single node
  - What will be partitioned? Tables, records, ....
- Scalability
  - Partitioning allows load balancing of data and its access
  - Works best for NO-SQL data stores, e.g., key-value data stores
    - each partition contains a range of records
- Performance
  - Throughput improvement by concurrent access of distinct partitions
    - requires proper allocation of records to partitions!
      - avoid hot spots that receive the bulk of the queries (skewed workload)
      - or combine with replication





### Partitioning of key-value data

Define a mapping  $\pi: K \to \{n \mid 0 \le n < N\}$  from the set of keys *K* to the set of node numbers. Note that  $|K| \gg N$ 

- by key range
  - for totally ordered keys
  - continuous ranges, boundaries need not be evenly spaced
  - efficient range queries
  - *"hot ranges"* due to skewed workload
- by hash of key
  - need not be a cryptographic key (one-way function)
  - achieves uniform distribution, evenly spaced boundaries
  - inefficient range queries
  - *"hot keys"* remain a problem

#### Assumes a 1-1 relation between partitions and nodes!





# **Rebalancing partitions**

#### Rebalancing may become necessary due to

- Increased number of queries (workload per node)
- Increased partition size
- Node failure

### Change N and (thereby) redefine the mapping $\pi$

- Redistributes records over partitions
  - especially  $\pi(k) = (hash(k) \mod N)$  is inefficient

### Allow multiple partitions per node.

- Introduce separate mappings: one for keys to partitions, and one for partitions to nodes
  - as a result partitions and nodes can be independently added or removed
- Split (merge) partitions that become too large (small)
- Load balancing by moving entire partitions between nodes



# **Quality drivers for replication**

Introducing replication into the architecture of a data store is mainly motivated by the following quality drivers

- Reliability
  - Replication removes a single point of failure and allows consensus protocols to deal with corrupted data (majority voting, etc.)
- Availability
  - If the probability that a single server becomes inaccessible is p, then the availability with n servers is  $(1 - p^n)$  (assuming server failures are independent, which need not be the case in, e.g., datacenters)
- Performance
  - Throughput improvement by concurrent access of distinct replicas
  - Latency reduction by accessing a nearby replica
  - Bandwidth reduction due to increased data proximity
- Scalability
  - Replication (mostly of compute resources) allows load balancing





## **Replication transparency**

- An architecture is *replication transparent* when the users of the system are unaware of the fact that several *replicas* (physical copies) of an object (resource) exist.
  - Refers primarily to data values; services that are implemented via multiple servers or multi-threaded servers are (by intention) distinguishable through increased performance.
  - Implies that clients identify only a single (logical) data object as the target of an operation and also expect only a single return value (as opposed to one for each replica on which the operation is performed).
  - Implies that the architecture is also *location transparent*, otherwise replicas could be identified by their location.
    - e.g., as revealed in their name.





## **Architectural concerns**

Replication (especially of data) is not for free: it raises additional concerns whose solutions incur costs:

- 1. What is the number of replicas and where should the replicas and replica managers be located?
  - Statically or dynamically resolved. In the latter case, also which party initiates the creation/destruction of replicas?
- 2. Whether and how to maintain consistency?
  - Ideally, all replicas have the same value (at least upon access)
  - Difficult because there is no notion of global time and state.
- 3. What architectural elements to use for storage and management?
  - Front ends, replica managers, caches, key-value stores, relational databases, load balancers, multicast infrastructure, logical clocks for time-stamping and versioning
- 4. What protocols to use for accessing (reading & writing) of replicas?
  - Push(server)-based or pull(client)-based, unicast versus multicast, group view, gossiping, degree of synchronization





# Basic model (CDK5, Wiesmann et al.)

- The basic model is multiple-client-single-server or multiple-client-multiple-server.
  - Servers represent a distributed data store with one or more access points for their clients
  - Operations invoked on the objects in the store are classified in two broad categories
    - Updates (also referred to as *write operations*) that *potentially* modify the state of the store
    - Queries, (also referred to as *read operations*) that inspect part of the state of the store.
  - Each server has a special entity, the *replica manager (RM)*, that is responsible for managing the local part of the data store.
    - Depending on the details of the architecture, this entity also serves as access point for clients. Alternatively, this task is performed by a separate front end (FE).





#### Figure 15.1 A basic architectural model for the management of replicated data



- RMs manage replicas of multiple objects, and each object is managed by a subset of all RMs. The size of the subset determines the number of replicas of the object. Often, each RM manages every object.
- RMs apply operations to their replicas as indivisible actions, and all actions can be recovered.
- The state of an RM is completely determined by its initial configuration and the sequence of actions performed.

Instructor's Guide for Coulouris, Dollimore and Kindberg Distributed Systems: Concepts and Design Edn. 4 © Pearson Education 2005

### **Basic architecture: behavioral view**

#### 1. Request phase

- Requests accepted by FE are communicated to a fixed RM (passive)
- Requests accepted by FE are multi-casted to all RMs (active)

#### 2. Coordination phase

- RMs determine whether, by which RMs and in which order (FIFO, causal, total) requests are performed; ordering can be enforced through delivery mechanism (usually total order)
- 3. Execution phase
  - RMs tentatively, i.e. undo is possible, execute the requested operation
- 4. Agreement phase
  - RMs reach consensus on the effect of the operation and commit
- 5. Response phase
  - Some RMs respond to the FE which in turn replies to the client

#### Beware: Phases need not be executed in this order





# Agenda

- Introduction
- Consistency models
  - Single-server paradigm
  - Conflicting operations
  - Data-centric models
  - Client-centric models
- Replica management
- Architectural case studies





# **Consistency models**

1. Are contracts between data store and clients

- Consistency is a property of the data store as a whole; for individual data items we speak about coherence for which there can be separate models.
- 2. Define the unit of consistency (the conit)
- 3. Determine the outcome of a sequence of read/write operations performed by one or more clients
  - Results obtained by individual clients
  - Resulting state of the store
- 4. Can be classified in two broad categories:
  - Data-centric models
  - Client-centric models





## Single-server paradigm

- The key idea behind consistency is that, for all parties involved, the operations appear as if they were performed as indivisible actions by a single server, i.e., in the same order and having the same effect.
  - *"the same effect"* requires that queries return the same value and updates leave the data store in the same state.
    - This state is a logical concept, because, in practice, it can occur that there is never a moment in time at which all replicas hold values in accordance with the state. However, if clients cease to access the store and no RM crashes, all replicas eventually must assume the same value.
  - *"the same order"* requires the existence of a global, system wide, notion of time, that is used to totally order the receipt of operation requests by the data store.
    - In practice, synchronization of local clocks can only approximate this ideal, which may or may not be sufficient (beware True Time in Google Spanner <sup>(i)</sup>)
- Since neither demand can be fulfilled, consistency models basically delineate how far the system may deviate from this ideal.





# **Conflicting operations**

- Two operations are *conflicting* when the outcome of executing them as a sequence of two atomic actions may *potentially* differ for the two possible execution orderings.
- Conflicting operations come in two flavors:
  - read-write conflicts (different values returned)
  - write-write conflicts (store left in different states)
- The actual values involved in write operations can be such that no conflict arises.
  - Example: Client1::Write (X, 42) Client2::Write (X, 42)
  - To limit the amount of state information, replication management protocols, in general, do not take values of operations into account, only their order.





### **Execution constraints**

Consistency puts constraints on the interleaving of operations allowed to the single server, and therefore, by the single server paradigm, also on the behavior, i.e., order of operations, occurring at the RMs of the data store.

- R1: The interleaved sequence of operations meets the specification of a (single) correct copy of the objects
  - i.e., produces the same results (return values and internal state)
  - maintains system invariants
- R2: The order of operations in the interleaving originating from a single client is consistent with the order in which that client issued them (program order).
- R3: The order of operations in the interleaving is consistent with the global (real-time) ordering of the operations.





### **Concurrency and interleavings**







24-Oct-19

## **Sequential consistency**

Actual executions that satisfy R1 and R2 are *sequentially consistent.* 

- A data store with replicated objects is *sequentially consistent* when *all* its executions are sequentially consistent.
- Sequential consistency allows swapping the order of pairs of subsequent operations originating from distinct clients to obtain a single server execution (even when they are conflicting)
  - Reflects possible transmission delays between clients and the data store and between data store AEs.
- This is the strongest form of consistency that can be enforced without loosing full benefits offered by concurrency.
  - Where concurrent events are those that have no causal dependency





### **Definition from TvS**

A data store is sequentially consistent when

The result of any execution is the same as if the read and write operations by all processes on the data store where executed in some sequential order and the operations of each individual process appear in this sequence in the order specified by its program.





Actual executions that satisfy R1 and R3 are called *linearizable* 

- A data store with replicated objects is *linearizable* when *all* its executions are linearizable.
- Linearizability allows swapping the order of pairs of subsequent operations originating from distinct clients provided they do not conflict, i.e., they refer to distinct objects, or they are both read operations.
- Swapping conflicting operations at a single server, in principle modifies the result (return value of reads or resulting state of logical object).









Single server execution (the only one allowed by R3)





Single server execution (the only one allowed by R3)

m: RM

Read (X)

Return 1

 $c_2$ : CL



# Single server execution (1 out-of 2 allowed by R3)







#### Single server execution (the only one satisfying R3)





# Single server execution 1 out-of 2 allowed by R3







Single server execution does not satisfy R3, shows swapping non-conflicts does not alter linearizability





Single server execution: does not satisfy R3, but does satisfy R2.



conflicting operations swapped (allowed)



Initially:  $X = \langle 0, 0 \rangle \land Y = \langle 0, 0 \rangle$ 

Can also happen without failure Just slow propagation

# Single server execution



Reordering reflects arrival time at replica manager  $m_2$ 



Single server execution 1-out-of 6 Reordering does not help









#### Single server execution



R.H. Mak

2IMN10-REP 33

## Other forms of consistency

### Causal consistency

- Operations from different clients that are causally related (as defined by Lamport's happens before relation) may not be swapped.
  - are seen by all clients in the causal order.
- Determined by means of time stamps and vector clocks.
- Enforced by the delivery mechanism in the coordination phase
- Weaker than sequential consistency.
- Eventual consistency
  - In the absence of further updates and system failures all replicas eventually have the same state.
  - Weaker than causal consistency.





m3: RM  $m_2: RM$ *c*<sub>3</sub>: *CL*  $c_2$ : CL $c_1$ : CL $m_1: RM$ Write (X, 1) Write (X, 1) Return Return Read (X) Return 1 := X+1 Write (X, 2) happens before Return Write (X, 2) Return 📕 Read (X) No causal Write (X, 1) ≯ Return 2 consistency Return Read (X) Return 1

Initially:  $X = \langle 0, 0, 0 \rangle$ 

**Beware**: RM  $m_2$  is superfluous.

Only present to show chain of

causal dependence





#### Single server linear execution



### **Client-centric models**

#### Monotonic read consistency

• If a process reads the value of a data item *x*, any successive read operation on *x* by that process will always return that same value or a more *recent* one.

#### • Monotonic write consistency

• A write operation on a data item *x* is *completed* before any successive write operation on *x* by the same process.

#### • Read your writes

• The effect of a write operation by a process on data item *x* will always be seen by a successive read operation on *x* by the same process.

#### • Writes follow reads

• A write operation by a process on a data item *x* following a previous read operation on *x* by the same process is guaranteed to take place on the same or a more recent value of *x* than was read.









Monotonic read violation at  $c_2$ 

or written by another global time stamps determine whether it

delayed propagation

Initially:  $S = \langle ab, ab, ab, ab \rangle$ 



Clients can only append, RMs have two options

- forward the operation
- forward its result

In either case, for this operation the mobile client can observe that forwarding has not been done, because it sees a prefix on the second read

### **Client-centric models**

### Monotonic read consistency

• If a process reads the value of a data item *x*, any successive read operation on *x* by that process will always return that same value or a more *recent* one.

#### Monotonic write consistency

• A write operation on a data item x is *completed* before any successive write operation on x by the same process.

### • Read your writes

• The effect of a write operation by a process on data item *x* will always be seen by a successive read operation on *x* by the same process.

### Writes follow reads

• A write operation by a process on a data item *x* following a previous read operation on *x* by the same process is guaranteed to take place on the same or a more recent value of *x* than was read.









#### Can be prevented by versioning updates





#### Can be prevented by making writes atomic

### **Client-centric models**

#### Monotonic read consistency

• If a process reads the value of a data item *x*, any successive read operation on *x* by that process will always return that same value or a more *recent* one.

#### • Monotonic write consistency

• A write operation on a data item *x* is *completed* before any successive write operation on *x* by the same process.

#### Read your writes

• The effect of a write operation by a process on data item *x* will always be seen by a successive read operation on *x* by the same process.

#### Writes follow reads

• A write operation by a process on a data item *x* following a previous read operation on *x* by the same process is guaranteed to take place on the same or a more recent value of *x* than was read







#### Initially (0,0)

#### X a private variable of $c_1$ $m_2$ : RM $c_1$ : CL $m_1: RM$ Write (X, 1) Return ← Client moves to another location, i.e., rebinds Read (X) Return 0 ← Write (X, 1) Return

Read-your-writes violation



#### Can be prevented by making writes atomic

### **Client-centric models**

#### Monotonic read consistency

• If a process reads the value of a data item *x*, any successive read operation on *x* by that process will always return that same value or a more *recent* one.

#### • Monotonic write consistency

• A write operation on a data item *x* is *completed* before any successive write operation on *x* by the same process.

#### • Read your writes

• The effect of a write operation by a process on data item *x* will always be seen by a successive read operation on *x* by the same process.

#### Writes follow reads

• A write operation by a process on a data item x following a previous read operation on x by the same process is guaranteed to take place on the same or a more recent value of x than was read.





Of course m is the "Ultimate Question of Life, the Universe and Everything"!!! See: Douglas Adams, The Hitchhikers Guide to the Universe

### Literature

- For the preparation of this slide set we have used material from various sources.
- Tanenbaum, van Steen: *Distributed Systems: Principles and Paradigms,* 
  - Chapter 7
- Coulouris, Dollimore, Kindberg, *Distributed Systems: Concepts and Design,*
  - Chapter 18 (basic model + gossip architecture)
- M. Wiesmann, F. Pedone, A. Schiper, B. Kemme, G. Alonso, <u>Understanding Replication in Databases and Distributed Systems</u>,
  - basic model
- R. Ladin, B Liskov, L. Shrira, and S. Ghemawat

Providing high availability using lazy replication

• gossiping



