

*The* UNIVERSITY of EDINBURGH  
SCHOOL *of* INFORMATICS

**CS4/MSc**

# **Distributed Systems**

Björn Franke

bfranke@inf.ed.ac.uk

Room 2414

(Lecture 14: Replication, 20th November 2006)

# Replication

Replication is a technique for **enhancing** services: multiple copies of data are maintained at multiple computers. This can lead to:

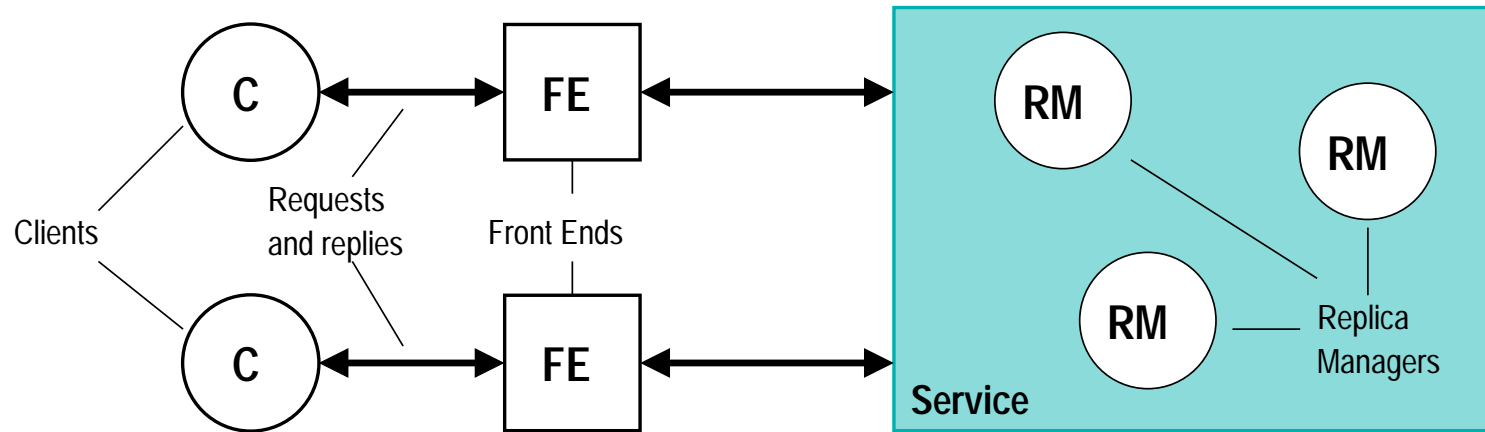
**Performance Enhancement:** For example, caching data in the web provides improved performance in terms of response time to users and lower utilisation of servers. In some cases the service itself may be replicated (e.g. DNS names may correspond to several IP addresses).

**Fault Tolerance:** A fault-tolerant service always guarantees strictly **correct** behaviour even in the presence of a certain number and type of faults. Correctness here may concern the integrity of the data with respect to client operations, and/or timeliness of the response.

**Increased Availability:** The availability of a service is the probability that a response is obtained within a reasonable time bound. Delays can be due to service features such as data locking as well as server and communication infrastructure failures.

## System Model

When data is replicated it should be done in a **transparent** manner: the client should not normally be aware that multiple *physical* copies of the data exist, either in terms of submitted requests or returned values. This is achieved by interposing a **front end** between the client and the service.



There are also issues of **consistency** between the replicas although the degree of inconsistency tolerated will depend on the service and the motivation for the replication.

## Executing a request on replicated objects

In general five phases are involved in the execution of a single request on a replicated object. The ordering of phases may vary.

**Initiation:** The front end issues the request to one or more replica managers (RMs).

**Coordination:** The RMs coordinate in preparation for executing the request consistently. They may agree on whether the request should be executed and the ordering of this request relative to others.

**Execution:** The RMs execute the request, possibly **tentatively**, i.e. in such a way that they can undo the effects later if necessary.

**Agreement:** The RMs reach consensus on the effects of the request (if any) that will be committed.

**Response:** One or more RMs respond to the front end. In some systems it is the responsibility of the front end to collect responses from a collection of RMs and select or synthesise a response for the client.

## Request ordering

Depending on the application different ordering semantics may be appropriate for the handling of requests. These are related to the possible orderings in multicast:

**FIFO ordering:** If the front end issues request  $r$  then request  $r'$  then any correct RM that handles  $r'$  handles  $r$  before it.

**Causal ordering:** If the issue of request  $r$  happened-before the issue of request  $r'$ , then any correct RM that handles  $r'$  handles  $r$  before it.

**Total ordering:** If a correct RM handles request  $r$  before request  $r'$ , then any correct RM that handles  $r'$  handles  $r$  before it.

Some applications require *hybrid* orderings, for example

**Forced ordering:** Requests are handled in an order that is both *causal* and *total*.

## Replication in CORBA

- An **object group** is a set of objects (usually instances of the same class) that process the same set of responses concurrently. Each returns a response but clients do not need to be aware of the replication.
- The client objects invoke operations on a single, local object which acts as a proxy, not for a single remote object, but for the group of objects.
- The proxy sends the invocation to each of the members of the group using group communication, and collects responses.
- Several systems support replication in CORBA in the form of object groups.
- For example, **Electra** is a CORBA-compliant object group system. Here proxies can be run in **non-transparent** or **transparent** mode depending on whether the client is replication-aware or not. The ORB is suitably extended to manage object groups.

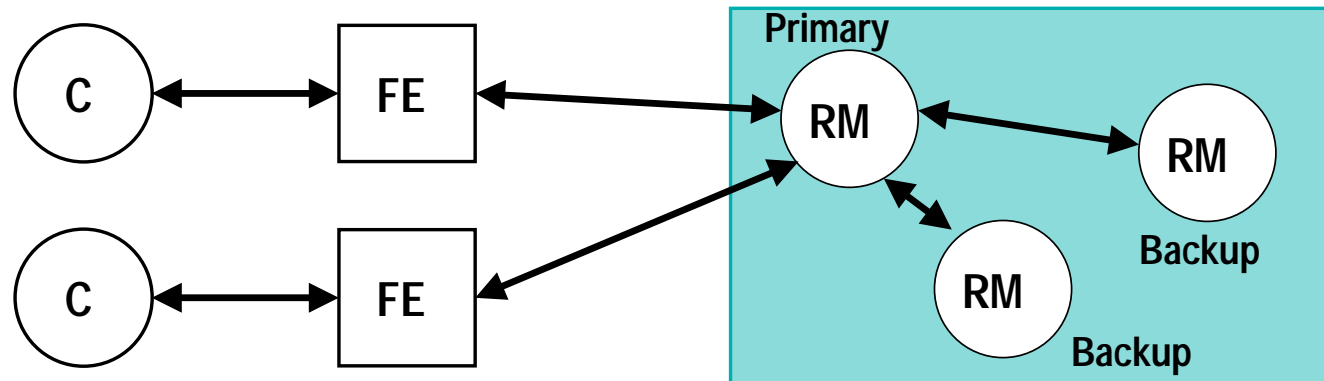
## Replication for Fault Tolerance

- A **fault tolerant** service based on replication should be able to keep responding despite failures and clients should not be able to tell the difference between the replicated service and one provided by a single correct RM.
- Care is needed to avoid anomalies with respect to the consistency of the data.
- Various notions of correctness have been defined such as:

**Linearizability:** for any execution there is an interleaving of the client operations that meets the application expectations for a single correct copy of the objects, and maintains a relative ordering of the operations which is consistent with the real times at which the operations occurred.

**Sequential Consistency:** is a weakened form of linearizability which imposes only that the relative ordering of operations is consistent with the program order in which each individual client executed them.

## Passive Replication (1)



In the **passive** model of replication for fault tolerance (also known as the **primary-backup** model), one RM is distinguished as the **primary** one. All front ends communicate with the primary RM which is responsible for executing all requests and then updating the other RMs, known as **backups** or **slaves**. If the primary RM fails it is replaced by one of the backups. At this point the surviving RMs must also agree on the set of operations that had been completed before the replacement primary takes over.



## Passive Replication (2)

The following sequence is followed if the primary is correct and guarantees linearizability:

1. *Request:* The front end issues a request to the primary RM. Each request contains a unique identifier.
2. *Coordination:* The primary atomically deals with each request in FIFO order; using the identifier it filters duplicate requests and resends the previous response.
3. *Execution:* The primary executes the request and stores the response.
4. *Agreement:* If the request is an update, the primary will propagate it, together with the response and the request id to all the backup RMs. The backups will send an acknowledgement to the primary.
5. *Response:* The primary responds to the front end; the front end responds to the client.

## Failure of the Primary

If the primary RM fails one of the backup RMs should take its place. The system will maintain linearizability if

- the primary is replaced by a unique backup;
- all the surviving RMs agree on which operations had been performed at the point at which replacement occurs.

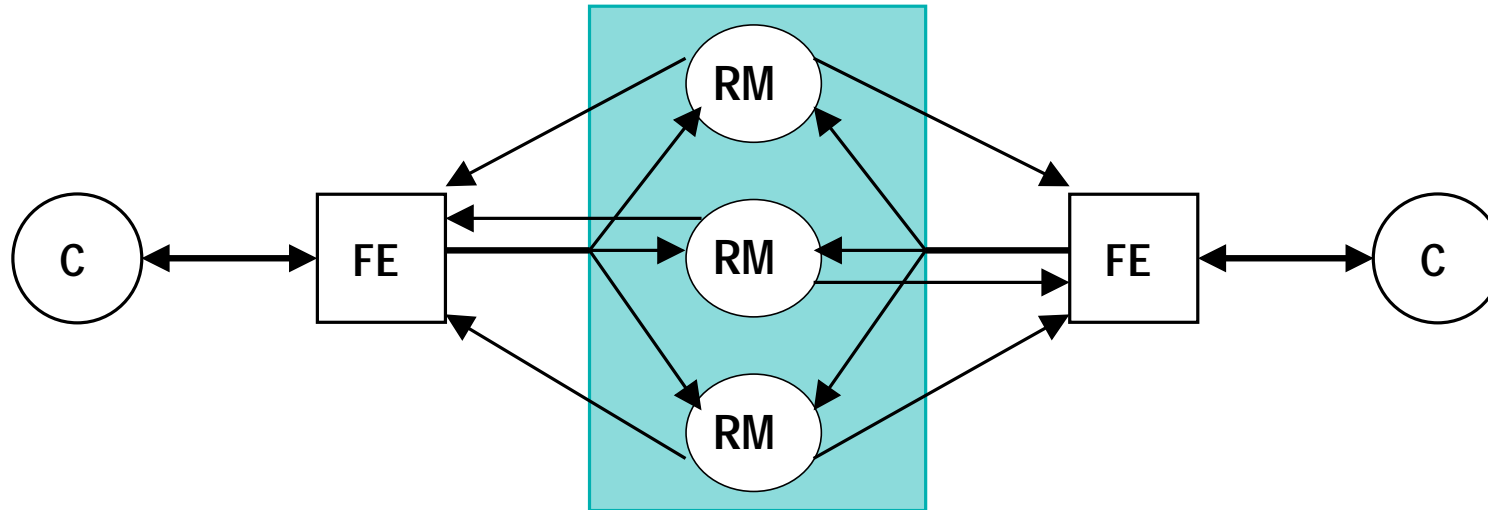
These requirements are met if the RMs are organised as a group and view-synchronous group communication is used to send updates to the backups:

- the view of the group will be consistent amongst the remaining RMs and will exclude the failed primary; a predefined function will be used to select the new primary from that view and this RM can assume the role.
- view-synchronous semantics guarantee that either all the backups or none of them will deliver any given update before delivering the new view.

## Passive Replication Issues

- Basing the updates on view-synchronous communication is costly in terms of overhead.
- A passive replication system cannot survive byzantine failures; to survive up to  $f$  crashes,  $f + 1$  replica managers must be included.
- Supporting fault tolerance in the front-end simply requires adding the functionality to allow the front-end to look up the new primary when the current primary does not respond.
- If clients are allowed to submit read requests to backup RMs congestion at the primary may be reduced but the linearizability property is lost. Instead the service will be sequentially consistent.
- A form of passive replication is used in the Sun Network Information Service (NIS). However replicated data is updated at a primary RM and propagated from there on a one-to-one basis to other RMs. This means that the service is not even guaranteed to be sequentially consistent.

## Active Replication (1)



In the **active** model of replication for fault tolerance, the RM are state machines that play equivalent roles and are organised as a group. Front ends multicast their requests to the group and each RM processes the request independently but identically and replies. Active replication can tolerate byzantine failures because the front end can collect and compare responses.

## Active Replication (2)

1. *Request:* The front end multicasts the request to all RMs, using totally ordered, reliable multicast, after attaching a unique identifier to it.
2. *Coordination:* The request is delivered to all correct RMs in the same order (by properties of the group communication used).
3. *Execution:* Every correct RM executes the request, producing the same result. This is guaranteed since the RMs are all state machines and they each handle requests in the same order. Each response contains the unique identifier.
4. *Agreement:* No agreement phase is needed, because of the multicast delivery semantics.
5. *Response:* Every RM sends its response to the front end. Differing policies can be enforced here. For example, the front end may respond to the client on the first response, discarding subsequent responses with the same identifier.

## Active Replication Issues

- This system is dependent on totally ordered and reliable multicast but when it is available the system achieves sequential consistency: all correct RMs process the same sequence of requests. The reliability ensures that all RMs process the same set of requests, and total ordering ensures that they satisfy them in the same order.
- If clients do not communicate with each other while waiting (e.g. if requests are synchronous and clients are not multi-threaded) then the requests will be processed in happened-before order.
- If clients are multi-threaded and exchange messages while waiting then to guarantee request processing in happened-before order we would have to replace the multicast with one that is both causally and totally ordered.
- The system does not achieve linearizability: the total order in which the RMs process request is not necessarily the same as the real time order in which they were submitted.

## Replication for High Availability

- When replication is included for high availability the intention is to give clients access to the service for as much of the time as possible.
- If data are replicated at two or more failure-independent servers then client software should be able to access data at an alternative server should the default server fail or become inaccessible.
- Furthermore replicated service can help to keep the response times low.
- Consistency issues are less of a priority; instead the priority is to return a response to the client. Thus collective agreement may not be reached before a response is made.
- Propagation of updates may be **lazy** rather than the **eager** style adopted in most fault tolerant replicated services.

# The Gossip Architecture

- The gossip architecture is a framework for providing a highly available service which replicates data close to the points where groups of clients request it.
- RMs exchange “**gossip**” messages periodically in order to convey the updates they have each received from clients.
- A gossip service provides two basic types of operation:
  - queries** are read-only operations; and
  - updates** modify but do not read the state.
- A front end can send a query or update to any RM that it chooses (based on availability and response time).
- RM may be temporarily unable to communicate with each other due to failures. Nevertheless the system guarantees a form of consistency called **relaxed consistency**.



## Gossip System Guarantees

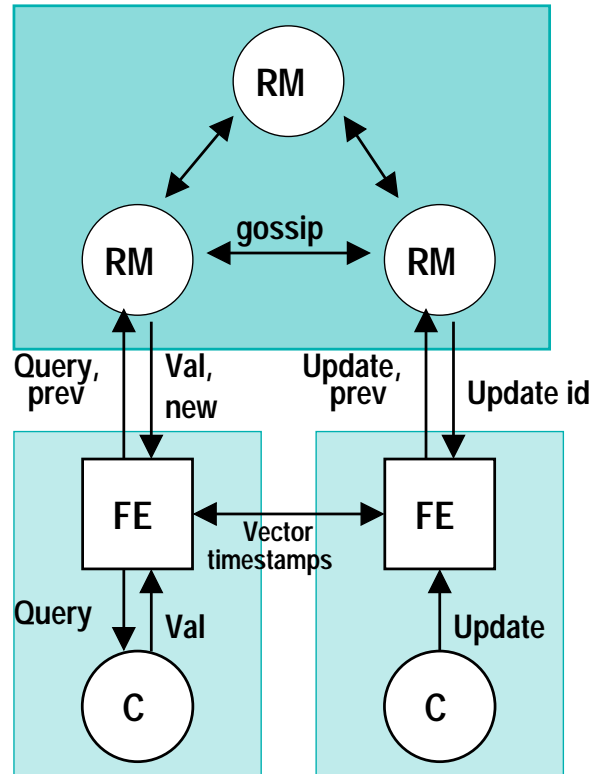
**Each client obtains a consistent service over time:** Responding to a query a RM provides a client with data that reflects at least the updates that client has observed so far.

Note that this is guaranteed even though a client may communicate with different RM at different times.

**Relaxed consistency between replicas:** All RM eventually receive all updates and they apply updates with ordering guarantees that make the replicas sufficiently similar to meet the needs of the application.

This is weaker than sequential consistency and may mean that clients observe stale data. It is supported by **causal update ordering**: if the issue of request  $r$  *happened-before* the issue of request  $r'$  then any correct RM handles  $r$  before  $r'$ . However stronger ordering conditions may be applied at higher operational costs.

# The Gossip Architecture



The gossip service front end handles client operations using an application-specific API and turns them into gossip operations (queries or updates). RM updates are **lazy** in the sense that gossip messages may be exchanged only occasionally. Each front end keeps a **vector timestamp** *prev*, reflecting the latest data values accessed by the client/front end. When clients communicate directly they piggyback their vector timestamps, which are then merged.

## Processing requests

1. *Request*: The front end sends the request to a RM. Queries are usually synchronous but updates are asynchronous: the client continues as soon as the request is passed to the front end and the front end propagates the request in background.
2. *Update response*: The RM replies to the front end as soon as it has received an update.
3. *Coordination*: The RM that receives a request does not process it until it can apply the request according to the required ordering constraints. This may involve receiving updates from other replica managers, in gossip messages.
4. *Execution*: The RM executes the request.
5. *Query response*: If the request is a query the RM responds at this point.
6. *Agreement*: The RMs only coordinate via gossip messages.

## Gossip Replica Manager: main components

**Value** This records the application state as maintained by this RM.

**Value timestamp** This represents the updates currently reflected in the value.

**Update log** All update operations are recorded in the log as soon as they arrive although they may be applied in a different order. An update is termed **stable** if it can be applied consistently with the desired ordering. An arriving update is stored until it is stable and can be applied. It then remains in the log until the RM receives confirmation that all other RMs have received the update.

**Replica timestamp** This timestamp represents those updates that have been accepted by the RM (not necessarily applied).

**Executed operation table** Maintained to prevent an update being applied twice and checked before an update is added to the log.

**Timestamp table** Contains a vector timestamp for each other RM, derived from gossip messages.

# Gossip Replica Manager

