



GINA CODY

SCHOOL OF ENGINEERING AND COMPUTER SCIENCE

Distributed System Design

COMP 6231 – Winter 2020

Concordia University

Department of Computer Science and Software Engineering

Instructor: R. Jayakumar

Distributed Event Management System(DEMS) – Group Project

Group members:

Sepehr Jalayer – 40126236 – jalayersepehr@gmail.com

Deep Shivam Gautam – 40138490 – deepshivamgautam@gmail.com

Seyed Pouria Zahraei – 40115175 – Zpouria@gmail.com

Table of Contents

Overview	3
Design Architecture	4
1.Data Consistency	5
2.Total Ordering	5
3.Dynamic Timeout.....	5
4.Fault tolerance (software bug)	5
5.UDP Multicast.....	5
6.UDP Reliability	5
7.Replica Recovery	6
Data Structures	7
Test Scenarios	8
Members tasks.....	9
Detailed Report (Sepehr Jalayer)	9
Detailed Report (Seyed Pouria Zahraei)	10
Detailed Report (Deep Shivam Gautam)	11

Overview

The distributed Event management system (DEMS) consists of three different servers which are located in different cities:

- Montreal(MTL)
- Sherbrooke(SHE)
- Quebec(QUE)

The clients of this system are of two types:

- eventManagers
- customers

We must ensure that these clients are connected to their own servers with Java RMI, and also the connection between our three servers are done through UDP/IP socket programming.

Manager specific functions:

- **addEvent()**: managers can only add events in their own server
- **removeEvent()**: managers can only remove events from their own server. *if an event was removed we must book another closest event for the customers registered in that event.
!!Needs UDP for server-server connection.
- **listEventAvailability()**: we must gather all events of a given type from all three servers.
!!Needs UDP for server-server connection.

Client/Manager functions:

- **bookEvent()**: customers can also book from other servers with a weekly 3 limit. !!Needs UDP for server-server connection.
- **getBookingSchedule()**: show the customers booking schedule.
- **cancelEvent()**: clients can remove an event from their own schedule. !!Needs UDP for server-server connection.
- **swapEvent()**: clients can swap a booked event with another event. (a bookEvent + cancelEvent) -> needs to be atomic

Clients are recognized with their ClientID (8 character): serverID (3char) + clientType(C/M) + 4 digit identifier.

Events are recognized with their eventType: Conferences/Seminars/Trade Shows + their eventID(10 character): serverID (3char) + eventSlot (M/A/E) + eventDate (DDMMYY).

Design Architecture

There will be three replica managers (RM) each containing different server implementations of our three servers and a database of that replica.

This document explains how high availability or fault tolerance distributed event management system is achieved. This system needs to fulfil the following criteria:

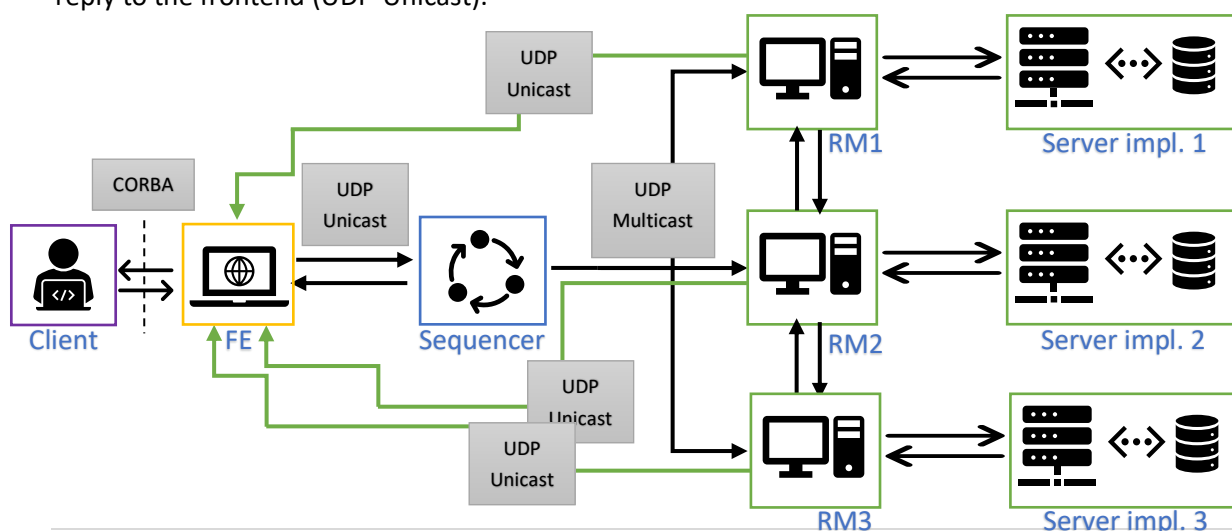
1. Data Consistency
2. Total Ordering
3. Dynamic Timeout
4. Fault tolerance (software bug)
5. UDP Multicast
6. UDP Reliability
7. Replica Recovery

Highly Available or Fault Tolerant Distributed Event Management System will be implemented over the CORBA project implemented as part of assignment2.

- We will implement a Front-End (FE) module which receives a client request as a CORBA invocation, forwards the request to the sequencer, receives the results from the replicas and sends a single correct result back to the client as soon as possible. The FE also informs all the RMs of a possibly failed replica that produced incorrect result.
- We will implement the replica manager (RM) which creates and initializes the actively replicated server system. The RM also implements the failure detection and recovery for the required type of failure.
- We will implement a failure-free sequencer which receives a client request from a FE, assigns a unique sequence number to the request and reliably multicast the request with the sequence number and FE information to all the three server replicas.

A simple request – response flow is mentioned below:

Front-End module would accept the Client Request (CORBA invocation). Front End will then send the request to the sequencer (UDP-Unicast) and the sequencer send the messages with a sequence id to replica managers (UDP-Multicast). It requires totally ordered and reliable multicast so that all RMs perform the same operations in the same order. Then RMs will process each request identically and send it to the corresponding servers. Specific servers are going to execute that request and going to reply to the frontend (UDP-Unicast).



1.Data Consistency

In order to achieve data consistency, we need to make sure that:

- Every request is executed by all or none of the replicas. (UDP Reliability, Dynamic Timeout)
- Every request is executed in the same order in all the replicas. (Total Ordering)

2.Total Ordering

Total ordering can be achieved using sequence numbers attached by the sequencer to every request. We use FIFO for the order of the processing.

Also, in case of a request lost, any RM receiving new request with any sequence number will check if the new sequence number is subsequent to the last executed request, if not; it will ask other RMs for the missing sequence numbers.

Other things that will help total ordering are:

- Dynamic timeout
- UDP Reliability
- UDP Multicasting

3.Dynamic Timeout

Each FE process initializes with a constant timeout (e.g.: 10 seconds). Then after execution of the first request it will change the timeout to the longest response time multiplied by 2.

If timeout reaches and FE does not get response from a server after 3 times, it will conclude that the server has crashed and informs other RMs to begin Recovery of that replica.

If timeout reaches and FE does not get response from any of the RMs, it will conclude that the request was lost and did not reach the RMs, therefore sends the request again(with the same sequence number).

4.Fault tolerance (software bug)

This is achieved through Front-End. The FE gathers the responses from all 3 RMs. It then decides the correctness simply by taking the majority. If an RM send 3 wrong responses, FE notifies all the RMs of the faulty RM in order to take the necessary actions.

Also each replica should be connected to a unique server implementation which includes all the methods and the data for Montreal, Quebec, Sherbrooke.

5.UDP Multicast

To decrease the number of requests send through the network, the connection between the sequencer and the RMs, also the connection between RMs are sent through UDP-Multicast.

6.UDP Reliability

UDP by itself is unreliable, therefore, we need to make the necessary steps to make it reliable.

- Timeout-resend: By the FE in order to overcome the request lost
- Multicasting each new request received by an RM to other RMs.
- If a new request received by an RM was not subsequent to the last executed request(there might have been a/some request loss), it will ask other RMs for the missing request if exists and then runs it/them before the new request.

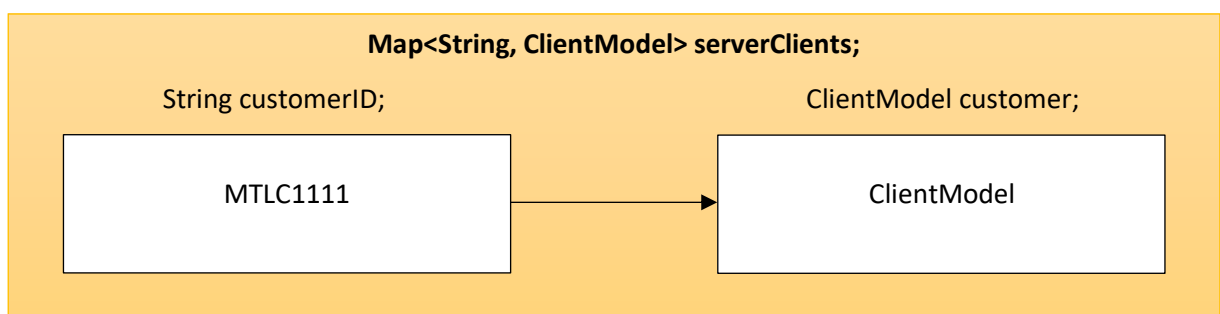
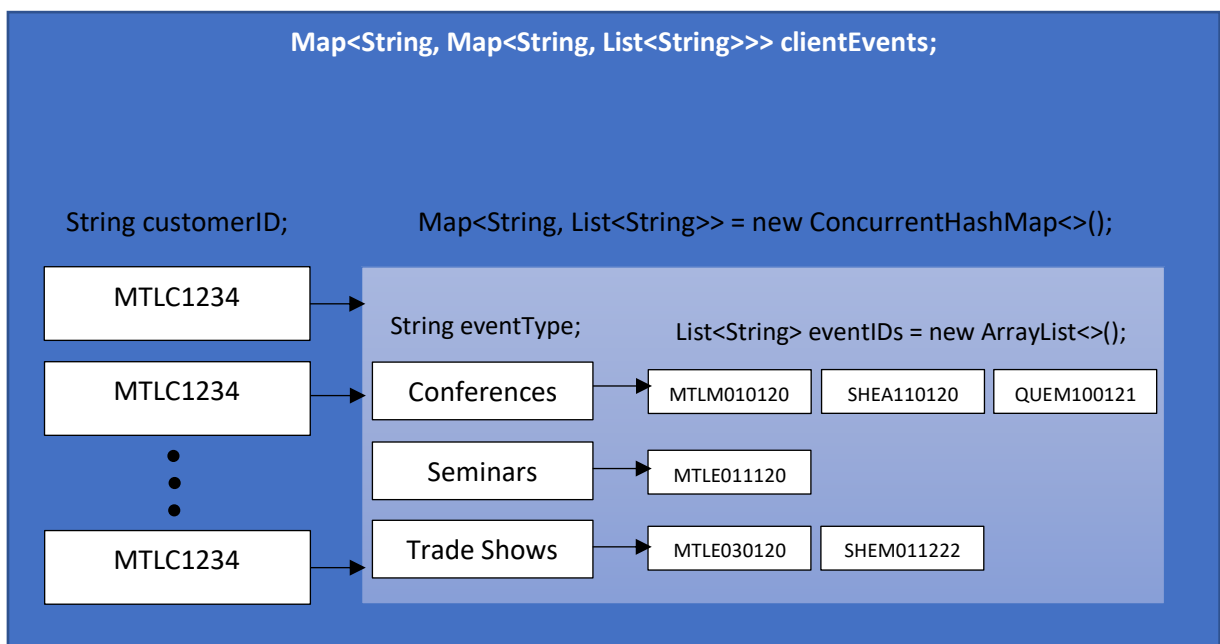
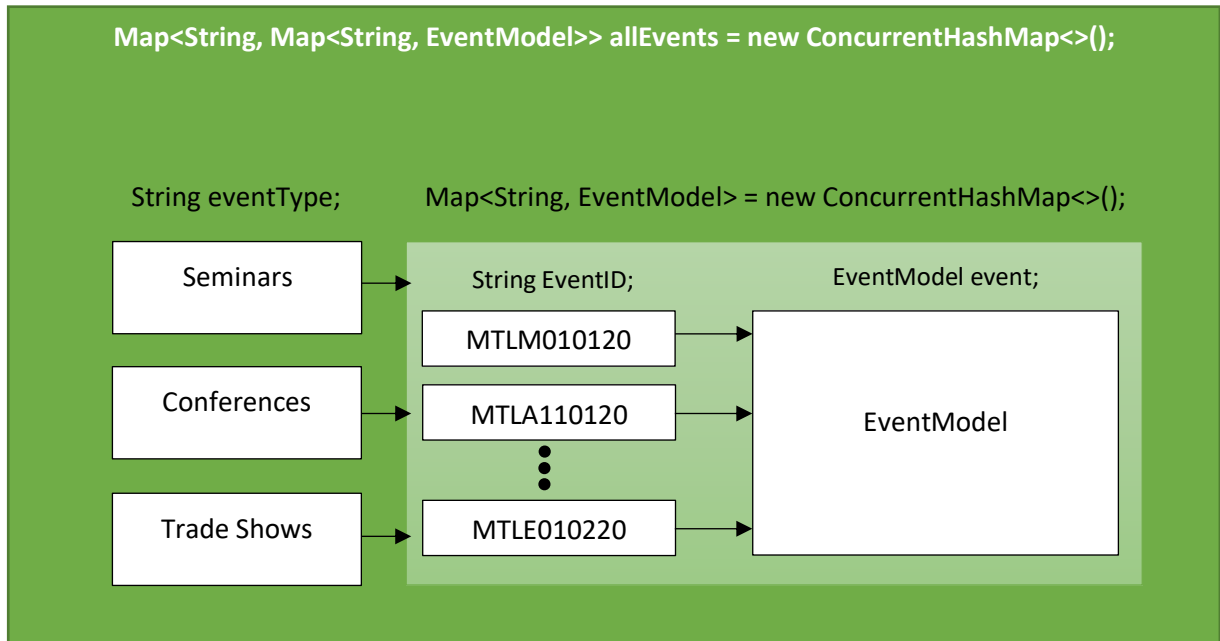
7.Replica Recovery

In order to achieve high-availability we need to ensure that a server crash can be tolerated. To do so, whenever a crash was detected (by FE) after not receiving a response from an RM for 3 times, or (by other RMs) after not receiving a heartbeat for some time (e.g.: 10 seconds) they will initiate a replica recovery for that RM.

In the replica recovery mode/state: the RM, asks other RMs for all the processed requests up to that point of recovery, then it will execute them one by one. After finishing the list of past requests, the replica is up and running.

Data Structures

All the data is maintained within each server, using three Map structures shown in the figure below.



Test Scenarios

#	Type of Test	Scenario	Cases
1	Login	UserName	1.Event Manager ID
2	Menu Items		2.Customer ID
3		Logout	1.Log out menu Item
4	Event Manager	addEvent()	1.invalid EventID -> not added 2.new EventID -> added 3.Existing EventID (LowerCapacity) -> not allowed 4.Existing EventID (HigherCapacity) -> capacity Updated 5.Duplicate Event -> not happening 6.EventID of Other Servers -> not allowed
5		removeEvent()	1.invalid EventID 2.EventID not exist 3.Event without anyone registered -> removed event 4.Event with someone registered -> Removed event + registered to same eventType if possible (UDP if needed) 5.EventID in other servers -> not allowed
6		listEventAvailability()	1.list all events of a given type from all three servers (UDP needed) 2.Event type is forced correctly with showing only options available
7		Ask for customerID	1.Access Customer methods
8	Event Manager + Customer	bookEvent()	1.on own server -> allowed 2.if event is full -> not allowed 3.on other servers -> only three in a week (UDP needed) 4. invalid EventID -> not allowed
9		getBookingSchedule()	1.Show booking schedule of customer 2.invalid customerID -> not allowed 3.customer not exist ->ok
10		cancelEvent()	1.cancel on own server -> ok 2.cancel on other server -> ok(UDP needed) 3.cancel a not registered event -> error shown 4.invalid eventide -> not allowed

11	swapEvent()	1.swap with a non-existent event -> not allowed 2.swap a not registered event -> not allowed 3.swap with a event from other server -> same conditions as a book event 4.swap with an event from other server in same week (weekly limit reached)
----	-------------	---

Members tasks

- Sepehr Jalayer: Front-End and Client and the server implementations in the replica 2
- Seyed Pouria Zahraei: Replica Managers and the server implementations in the replica 1
- Deep Shivam Gautam: Sequencer and the server implementations in the replica 3

Detailed Report (Sepehr Jalayer)

FE (Front End)

In this project we used a front-end as a median to communicate with the Client, RMs and the sequencer. The front-end uses CORBA interface to send/receive message to/from client. For the system to become fault tolerance or highly available, the FE calculates the majority of responses and detects every crash/bug due to response from each RM.

What was done?

1. Basically the FE serializes the request from client in a message format which was standardized by the
`group(Sequenceid;FrontIpAddress;MessageType;function(addEvent,...);userID;newEventID;newEventType; oldEventID; oldEventType;bookingCapacity)` to be sent to the sequencer and then to the RM.
2. Designed a dynamic timeout for the FE to wait for the response from RMs, depending on the longest wait time.
3. Created two classes in the Front-end package to creating an object for the request which can then be serialized for UDP. Another class for parsing the String message sent from Rms as a response and put it in an Object for convenience.
4. Wrote the majority detection algorithm which works on the parsed responses from each RM to detect the bug and the correct output to show to client.
5. Created the CommonOutput.java class for simple standardization of the responses from all the implementations.
6. Standardized a message type which is a code for the communication between FE-RM / SQ-RM to capture different scenarios.
7. Since UDP is unreliable, I created a retry condition if not getting a response from any RM after the timeout.
8. Included the Front-End Ip in the message format so that RM can send the responses to this IP via UDP Unicast.

How it works?

The Fe serializes the requests from the client into a known format and sends it to the sequencer via UDP unicast.

Then waits for the sequencer to return the sequence ID for that request and stores it for later to be able to assign RM responses to the corresponding requests.

Then waits for a dynamic timeout (initially 10 sec). *if the timeout reaches, any RM that did not send a response will then be counted as timed out (if the counter of each RM reaches 3 we detect a crash in that RM)

When we received all the responses and parsed them through a designed object (RmResponse.java) then we decide the majority. *each Rm response that don't match others will be counted as a bug in that RM (if the counter reaches 3 we detect that RM has software bug)

Both the bug and crash will be informed to all the RM via UDP multicast to the multicast ip of 230.1.1.10 and a corresponding 2 digit message type (e.g: 23-> RM3 has crashed or 13->RM3 has bug).

Why this approach was chosen?

For the Client-FE corba seems convenient because by using it we can use clients with different programming languages.

Also, putting all the standard communications and message formats in class which every module can have access to kept the code reworks needed to the minimum since everyone had to implement the same class which strictly forced them to use the standard messages.

Detailed Report (Seyed Pouria Zahraei)

RM

In this project we used 3 replica managers which, is capable of providing fault tolerance or highly availability by selecting the feature when the server system is initialized.

What was done?

We have 3 RMs, which are running in different machines and each RM is connected to 3 different servers using RMI connection. The connection between RMs and Sequencer is using multicast receiver and between each RM and FE is with UDP unicast. Between each RM also the communication is happening using Multicast receiver.

How it works?

Each RM once they receive a message, they will multicast the message to other RMs and will add the message to a queue and a HashMap which the key of this HashMap is sequenceID and the value would be the message.

We have a separate thread where is handling the execution of messages using priority queue.

There is connection between RMs and 3 servers where each RM sends request to servers using RMI connection.

We Created a class for having a standard format for messages between RM and FE.

Message format:

*Sequence id;FrontIpAddress;Message Type; function(addEvent,...);userID;
newEventID;newEventType; oldEventID; oldEventType;bookingCapacity*

Using Message Type in RM we can indicate the type of request(00- Simple message,...).

Each RM in case there a bug in any server, will inform which RM has bug.

In case the servers of each RM have been crashed, the RM will restart each server and will reload the databases in each server.

In case of each RM has lost a message, they will ask other RMs to update their HashMap, using the sequenceID, means they ask only the missing sequenceID since they are unique and then replica will start to execute the requests.

Why this approach was chosen?

By using RMI for sending requests to server it would be much simpler to compare with UDP. And also, by multicasting each request to other RMs we make sure that other RMs will not miss the request and also in order to achieve high-availability in case of one of servers are down we can restart the servers and reload the servers using recovery HashMap.

Detailed Report (Deep Shivam Gautam)

Sequencer

Sequencer is a middle ware that runs on unicast to get the information from Front-End and uses Multi-cast to send the data to 3 RM's.

What was done?

Sequencer is equipped with a datagram Socket at port 1333 and INET Address from FE. The received message is split and put in to usable parts then one part of the message is sent back to FE with the same address from which Sequencer received the data and other part along with IP and incremented sequencer id is multi-casted to all 3 RMS via INET address 230.1.1.10

How it works?

The data sent via a FE is captured and split into parts using semicolon. The first part is Sequencer id along with the rest of the message. This sequencer id is sent back to FE where it verifies number of attempts made to send the data on the other hand the sequencer id is incremented and added to remaining message along with the IP of FE which is sent to the RM.

Why this approach was chosen?

UDP Multicast is used at this stage to reduce the number of request sent through the network.

To achieve total ordering we have added sequence numbers to every request.
Therefore it can be identified at RM and operated as required.

In case of loss of request or any issues with request we will increment the sequencer id and send it back to Front end, so the data can be resent.