



Università
Ca' Foscari
Venezia

Scalability and Performance Analysis of a Web Application

Laureanda/Laureando

Federico Carraro, 847386

Andrea Brion, 860595

Anno Accademico

2024/2025



Università Ca'Foscari Venezia

2:

introduzione e scopo del progetto

Step 1

Create a web application that allows for searching a movie in the database and provides the information on the movies, the main actors and directors and the average rating, length etc. Design your application with the tools that you prefer.

Step 2

Create a query set. Let us assume that the probability that a film is searched is proportional to the number of rating that it has received.

Create a list of 10,000 queries of movie titles (where entries can be duplicated) sampled according to this rules.

The goal of the project is to identify the element in the network that limits its performance, the so called bottleneck.

Afterwards we'll define the optimal parameters and show how targeted resource increases can reduce the bottleneck or even shift it to another component.

Step 3

Perform a load test to assess the scalability of your web app in closed-loop.

How many users can it handle?

Use the queries that you previously created as input of your load test.

Step 4

Design the system in such a way that its scalability increases.

To this aim, identify the bottleneck in your implementation (either experimentally or from its model).

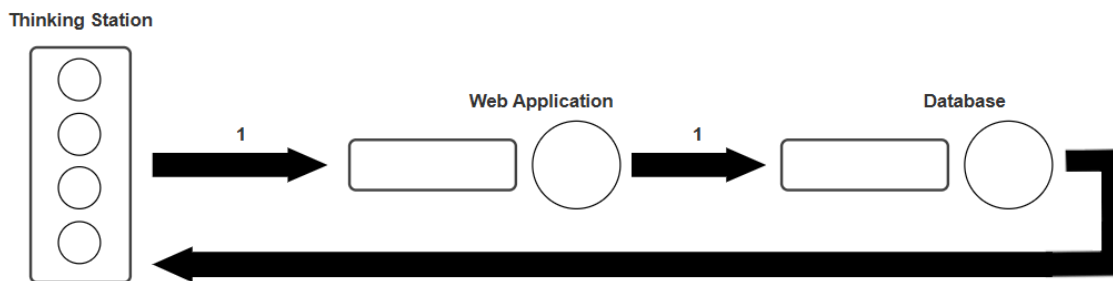
Then, test your designed system using JMT to study the expected response time as function of the number of users. What is the optimal number of users?



Università
Ca' Foscari
Venezia

6: *Analytical model:*

In order to perform load testing we decided to structure the system under test in an **interactive single-class closed-system** depicted in the following way:



Being a closed system queueing network the web server is going to serve a fixed number of users placed in the **thinking station** in which they send parallel and independent requests to the **web application station** that processes them.

The **database station** produces a response for each one and send back to each related user which will consult it for a certain amount of time called thinking time before doing another request.

We assume for the sake of simplicity that the web application and the database stations are single server with constant speed while users are statistically identical and so single-class. In addition, each station is always visited by every user so the relative visit ratio for each one is equal to 1.

3: *System architecture:*

Our web application is made up of a JavaScript file, given that our goal was solely to test performance, the back-end logic was implemented within this file: the application was designed to send 10,000 queries to the database.

Node.js enabled the JavaScript file to interact with the database, which was set up locally using PostgreSQL.

In accordance with the project guidelines, the data was sourced from the IMDB website.



Università
Ca'Foscari
Venezia

4:Load testing setup:

In order to perform load testing in our application we relied on Apache Jmeter load tester software.

We created a test plan, which defines the overall settings for the test, then added to it two distinct thread groups:

- one to simulate server-to-database requests
- one to simulate user-to-server requests

Now, we'll take a closer look at both.

The former one simulates the interaction between the web application and the database and it contains two main components:

- **Configuration element**, that establishes a connection to the database
- **JDBC sampler**, in which we've defined a query the web application use to retrieve data from the database and show it in the web application.

Worth mention that we've decided to set a maximum number of connections (into the configuration element) of 50, this because if the default value (1) was kept we've seen a error rate too high regarding the requests.

The latter simulates HTTP interactions between users and the server and it contains two main components:

- **HTTP sampler**, which allows sending HTTP/HTTPS requests to a web server and also parse embedded resources (e.g., images) when instructed, simulating realistic HTML content retrieval.
- **Flow Control Action sampler** configured for Pause only to add a **Constant Throughput Timer** that is the key factor in evaluating system behavior under a controlled load. This timer regulates the number of requests sent per minute or second (it dynamically adjust pacing to achieve a defined throughput), here we've set a random delay with an average and deviation to achieve a more realistic user interactions



Università
Ca' Foscari
Venezia

Right below is possible to see all the settings of the tests plan above described:

The screenshot shows the 'Thread Group' configuration window. The 'Name' field is set to 'Thread Group'. The 'Comments' field is empty. Under 'Action to be taken after a Sampler error', the 'Continue' radio button is selected. The 'Thread Properties' section includes: 'Number of Threads (users)' set to 1, 'Ramp-up period (seconds)' set to 1, 'Loop Count' with 'Infinite' unchecked and '10000' entered, and three unchecked checkboxes: 'Same user on each iteration', 'Delay Thread creation until needed', and 'Specify Thread lifetime'. The 'Duration (seconds)' and 'Startup delay (seconds)' fields are empty.

The thread group used for the server-to-database requests.
Here we've only set the number of queries that a user sends: 10000, as indicated in the project specification.

The screenshot shows the 'HTTP Request' configuration window. The 'Name' field is set to 'HTTP Request'. The 'Comments' field contains the text 'Qui decidiamo qual'è l'app web di destinazione'. The 'Basic' tab is selected. Under 'Web Server', 'Protocol (http)' is set to 'http', 'Server Name or IP' is set to 'localhost', and 'Port Number' is set to '3000'. Under 'HTTP Request', the method is set to 'GET' and the 'Path' is set to '/test'. The 'Content encoding' field is empty. At the bottom, the 'Parameters' tab is selected, showing a table with columns 'Name', 'Value', 'URL Encode?', and 'Content-Type'. The table is currently empty.

The JDBC configuration element, where it's possible to see the database's connection settings.



Università Ca' Foscari Venezia

We've used a local database in the configuration, in fact the URL specified is a local one

JDBC Request

Name:

Comments:

Variable Name Bound to Pool:

Variable Name of Pool declared in JDBC Connection Configuration:

SQL Query

Query Type:

Query:

```
1 SELECT DISTINCT
2     nb."author_name",
3     nb."author_firstname",
4     tb."original_title",
5     tb."isbn",
6     tb."language"
7 FROM   nb, tb
8 JOIN LATERAL unnest(string_to_array(nb."known_titles", ' ')) AS known_title("known") ON true
9 JOIN "titleBase" tb ON tb."known" = known_title."known"
10 WHERE nb."language" IN (
11     'English',
12     'French',
13     'German',
14     'Greek',
15     'Hebrew',
16     'Italian',
17     'Latin',
18     'Portuguese',
19     'Russian',
20     'Spanish'
21 )
```

Parameter values:

Parameter types:

Variable names:

Result variable name:

Query timeout (s):

Limit ResultSet:

Handle ResultSet:

The JDBC sampler.



Università Ca' Foscari Venezia

Thread Group

Name: Thread Group

Comments: Dato che vogliamo simulare un closed network

Action to be taken after a Sampler error

☒ Continue ☐ Start Next Thread Loop ☐ Stop Thread ☐ Stop Test ☐ Stop Test Now

Thread Properties

Number of Threads (users): 1

Ramp-up period (seconds): 1

Loop Count: ☒ Infinite

☐ Same user on each iteration

☐ Delay Thread creation until needed

☐ Specify Thread lifetime

Duration (seconds): 1

Startup delay (seconds): 0

The thread group used for the user-to-server requests.
We've put the flag to "Loop count: Infinite" to simulate a real situation where the users are not limited in the number of requests they can make.

HTTP Request

Name: HTTP Request

Comments: Qui decidiamo qual'è l'app web di destinazione

Basic Advanced

Web Server

Protocol (http): http Server Name or IP: localhost Port Number: 3000

HTTP Request

GET Path: /test Content encoding:

☐ Redirect Automatically ☒ Follow Redirects ☒ Use KeepAlive ☐ Use multipart/form-data ☐ Browser-compatible headers

Parameters: Body Data Files Upload

Send Parameters With the Request:

Name	Value	URL Encode?	Content-Type
------	-------	-------------	--------------

Here we see the HTTP Sampler, where we specified the protocol, the server name, the port (of the web application server), the type of request and the path to the resource that allows to the users to get the data.



Università Ca' Foscari Venezia

Flow Control Action

Name:

Comments:

Logical Action on Thread

☒ Pause ☐ Start Next Thread Loop ☐ Break Current Loop

Duration (milliseconds):

☐ Go to next iteration of Current Loop

Logical Action on Thread/Test

☐ Stop ☐ Stop Now

Target:

Here we see Flow Control Action.

Uniform Random Timer

Name:

Comments:

Thread Delay Properties

Random Delay Maximum (in milliseconds):

Constant Delay Offset (in milliseconds):

Here we see the Constant Throughput Timer, where we've set the two delays.



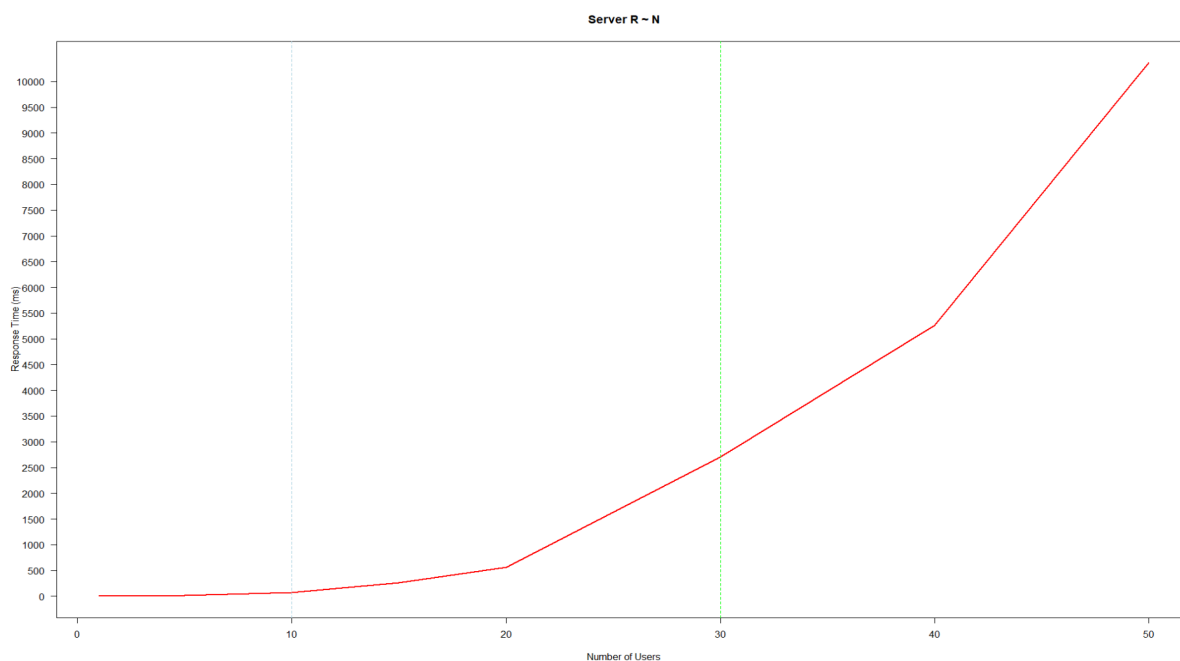
Università
Ca' Foscari
Venezia

5:

The results of **jMeter** obtained with the increase of the number of customer.
To better readability we've used rStudio to make the plots.
Analyzing the results we've found that

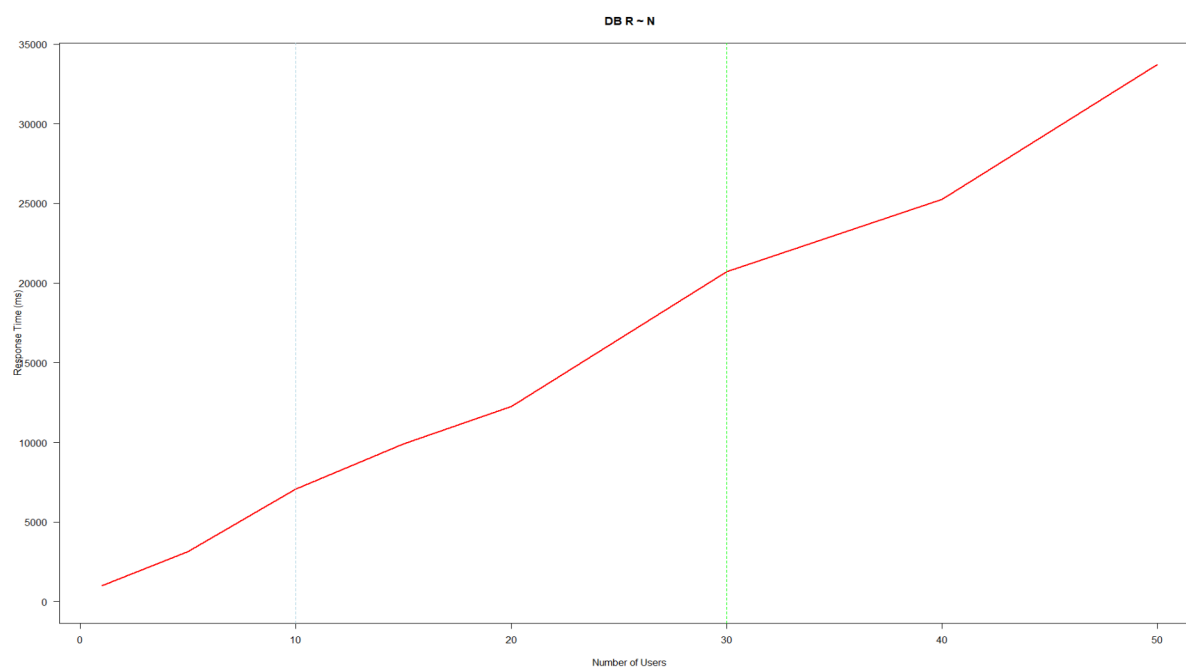
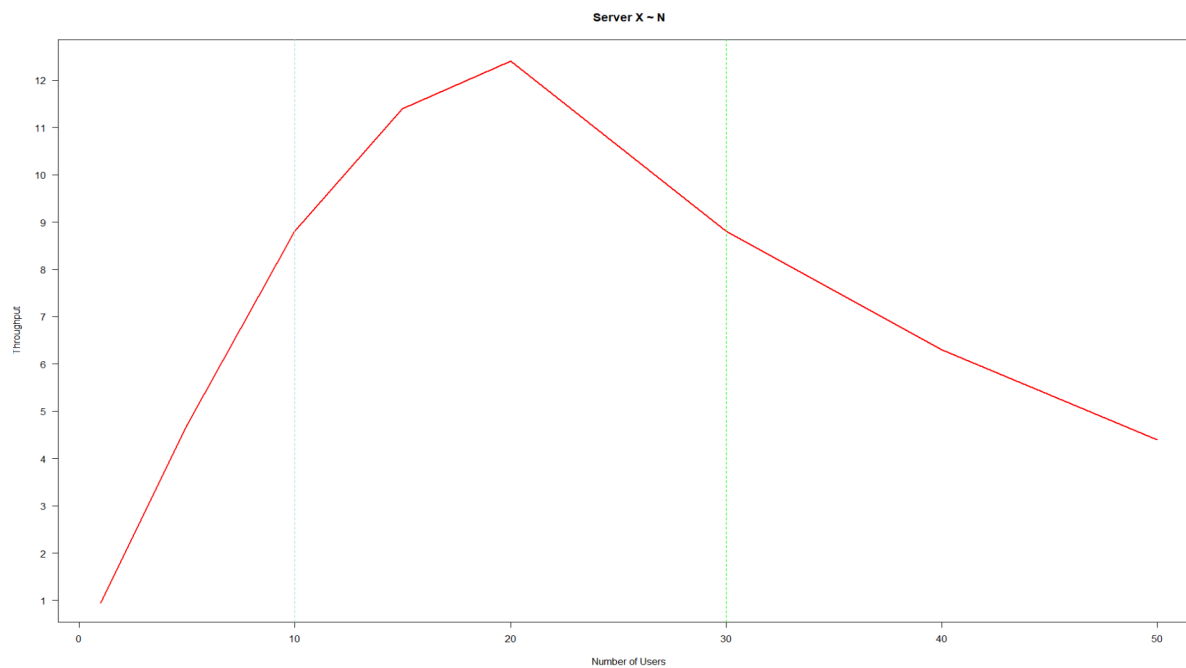
- the low load is with a number of customer below 10
- the mid load is with a number of customer below 30
- over 50 customer we have heavy load

As indicated by the dashed lines in the graphs that follows:



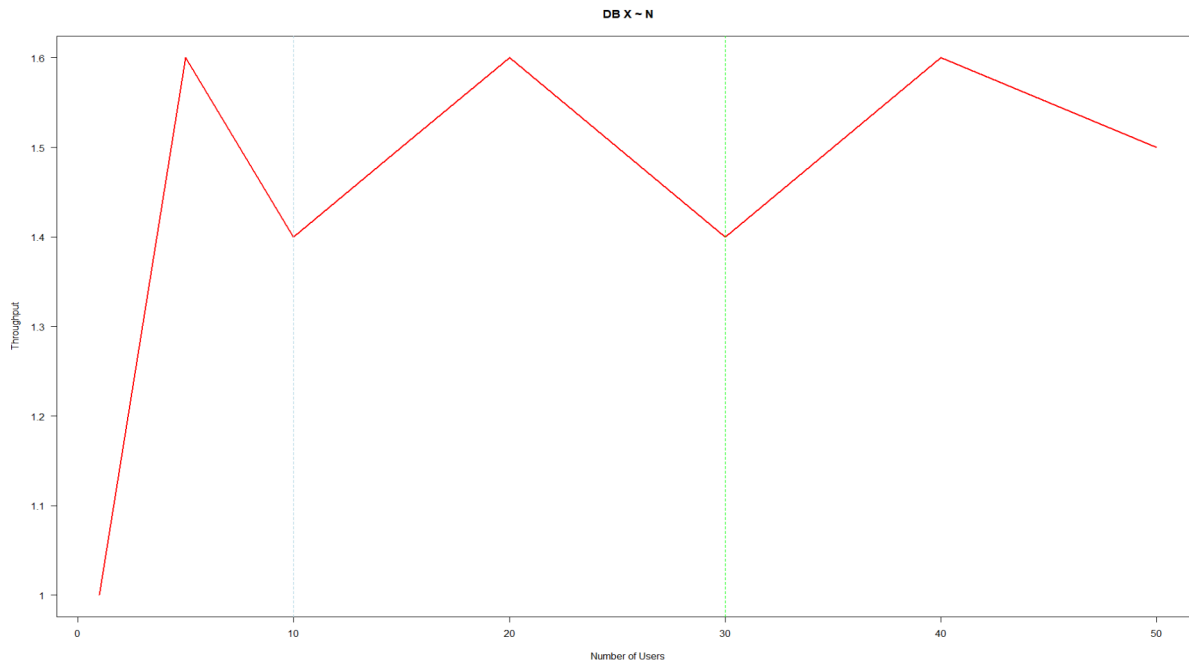


Università Ca' Foscari Venezia





Università
Ca' Foscari
Venezia



6.1: Model analysis by JMT

The settings in JMT touch various sections, let's dive into each of them:

- Classes, here we've set the network type and the types/classes of the requests, respectively: closed, single-class
- Stations, here the number and type of stations are set, respectively: thinking station, web-application station, database station
- Service Times, here the results obtained with jMeter are used to set each station service time
- Visits, here the visit ratios are set as one visit per station each visit to the reference station
- Reference station, here which stations is used as reference station: the thinking station
- What-if, here we set the bound of the experiment and the parameter to increase, respectively: 1 to 50 customer, number of customer

Right below the settings we've used described above:



Università Ca' Foscari Venezia

Classes Stations Service Times Visits Reference Station What-if Comment

Classes characteristics
Number, customized name, type of classes and number of customers (closed class) or arrival rate (open class). Add classes one by one or define total number at once. Higher number means higher priority

Number: 1

*	Name	Type	Priority	No. of Customers	Arrival Rate (λ)
1	Class1	closed	0	1	

The classes.

Classes Stations Service Times Visits Reference Station What-if Comment

Stations characteristics
Number, customized name and type of stations. Add stations one by one or define the total number at once. Load Dependent stations necessarily require the use of MVA.

Number: 3

*	Name	Type
1	Station T	Delay (Infinite Server)
2	Station DB	Load Independent
3	Station AW	Load Independent

The stations.

Classes Stations Service Times Visits Reference Station What-if Comment

Service Times
Input service times of each station for each class.
If the station is "Load Dependent" you can set the service times for each number of customers by double-click on "LD Settings..." button.
Press "Service Demands" button to enter service demands instead of service times and visits.
MULTICLASS MODELS: when for a station the per-class service times are different, the results are correct ONLY IF its scheduling discipline is assumed Processor Sharing (PS) and not FCFS (See BCMP Theorem).

*	Class1
Station T	1.0500
Station DB	0.9830
Station AW	0.0090

The service times.



Università Ca' Foscari Venezia

Classes	Stations	Service Times	Visits	Reference Station	What-if	Comment								
Visits Average number of visits to each station per class.														
			<table border="1"><thead><tr><th>*</th><th>Class1</th></tr></thead><tbody><tr><td>Station T</td><td>1.0000</td></tr><tr><td>Station DB</td><td>1.0000</td></tr><tr><td>Station AW</td><td>1.0000</td></tr></tbody></table>	*	Class1	Station T	1.0000	Station DB	1.0000	Station AW	1.0000			
*	Class1													
Station T	1.0000													
Station DB	1.0000													
Station AW	1.0000													

The visit ratio.

Classes	Stations	Service Times	Visits	Reference Station	What-if	Comment				
Reference Station The station is used to compute the system throughput and the system response time for each closed class . Performance metrics of open classes are always computed with respect to the arrival process . Visits at the Reference station can not be Zero. WARNING: the reference station for all closed classes is forced to be the same station.										
		<table border="1"><thead><tr><th>Class</th><th>Station</th></tr></thead><tbody><tr><td>Class1</td><td>Station T</td></tr></tbody></table>	Class	Station	Class1	Station T				
Class	Station									
Class1	Station T									

The reference station

Classes	Stations	Service Times	Visits	Reference Station	What-if	Comment
What-if analysis Solve models with increasing (or decreasing) number of customers of selected closed class.						
					Control Parameter:	Number of Customers
					Class:	Class1
					From (N1):	1
					To (N2):	50
					Steps (n. of executions):	10

The What-If.



Università
Ca' Foscari
Venezia

grafici di throughput, response time ed utilization
low, moderate ed heavy load

6.2: Analisi bottleneck mediante traffic equations, forced flow law e bottleneck law e grafici

Bottleneck

In order to spot the bottleneck we determine the average service demand for each station that is the total amount of service required by a customer to each station i for each visit it does at the reference station, assuming single server systems and independent service time distribution from the number of visits performed by a job at queue:

$$E[D_i] = E[V_i] * 1/\mu_i$$

and then the highest one determines the bottleneck station.

In our system we have that:

each station i is visited with visit-ratio: $E[V_i] = 1, \quad i = 1, 2, 3$

service demand web application: $E[D_{app}] = E[V_{app}] * 1/\mu_{app} = 1 * 0.0090 = 0.0090$

service demand database: $E[D_{DB}] = E[V_{DB}] * 1/\mu_{DB} = 1 * 0.9830 = 0.9830$

Therefore the database station is the bottleneck!

8: Calcolo optimal number of users sia analitico che grafico mettendo in risalto il throughput e response time con i relativi bound

Optimal number of users:



Università Ca' Foscari Venezia

After having identified the bottleneck, what happens to our system when the number of customers increases?

is going to be the first to have its utilization equal to 1, because the jobs tend to accumulate at that station and as a consequence limiting the system performances, namely the throughput.

In fact, the forced flow law states that the system throughput X_{sys} which corresponds to the reference station one is equal to $\mu_{\text{bottleneck}}/E[V]_{\text{bottleneck}}$

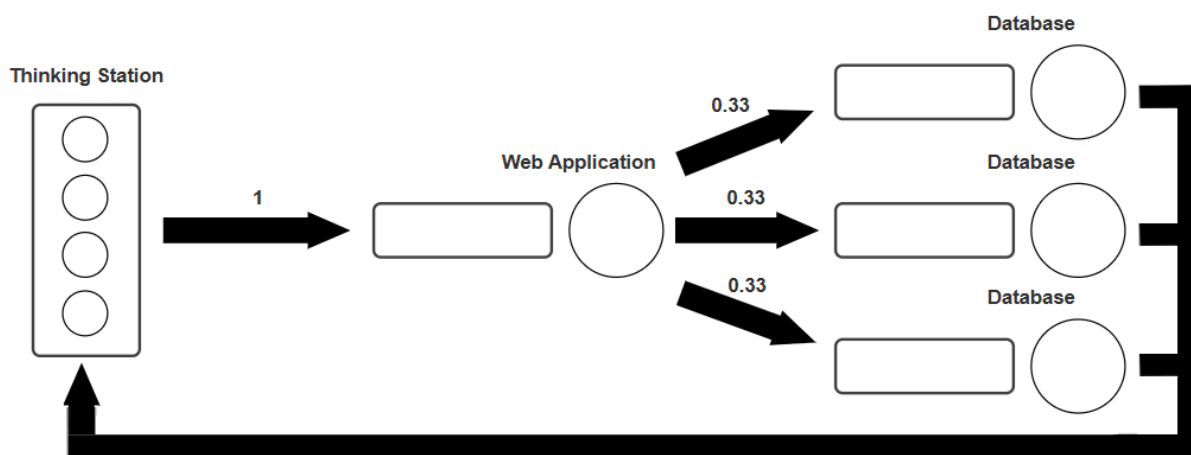
which imposes asymptotic limits to our system

9:ottimizzazione bottleneck via JMT con relative analisi throughput,response time e optimal number of users.

Following the literature, we have attempted (and succeeded) in reducing the bottleneck effects by increasing the number of stations of the same type as the bottleneck station.

We kept all settings the same for the additional stations, except for the visit ratio, which was set to 0.33 for each database station.

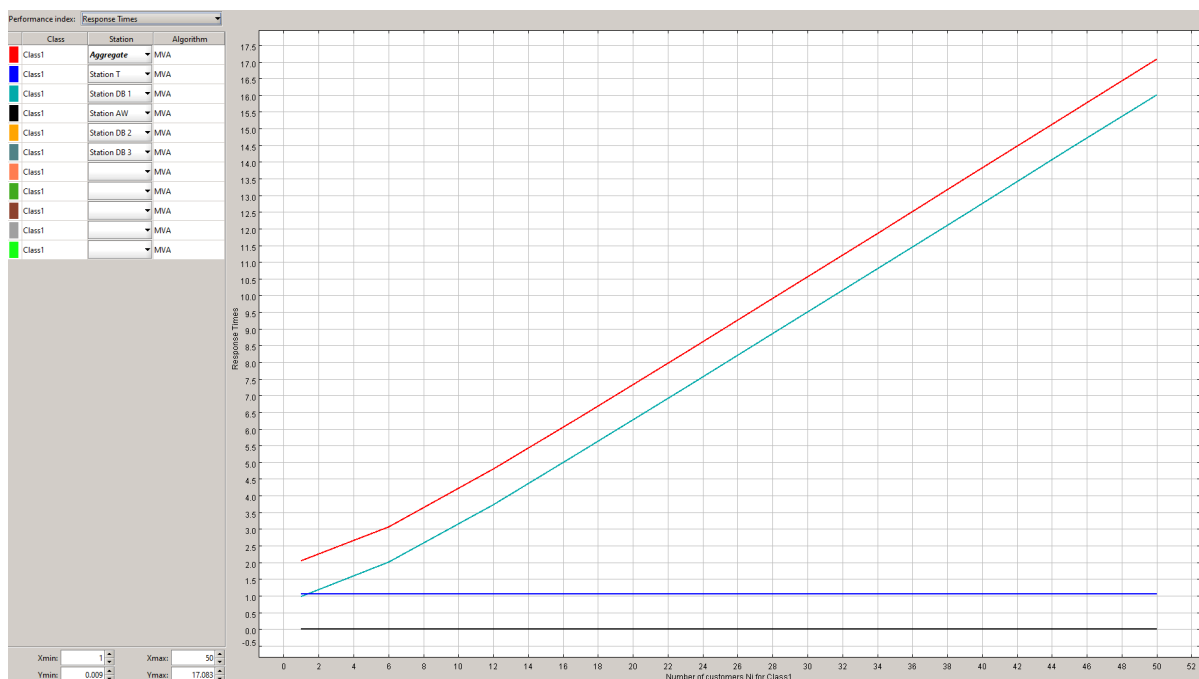
This value is derived by assigning an equal probability for a request to be routed to one of the three available stations.





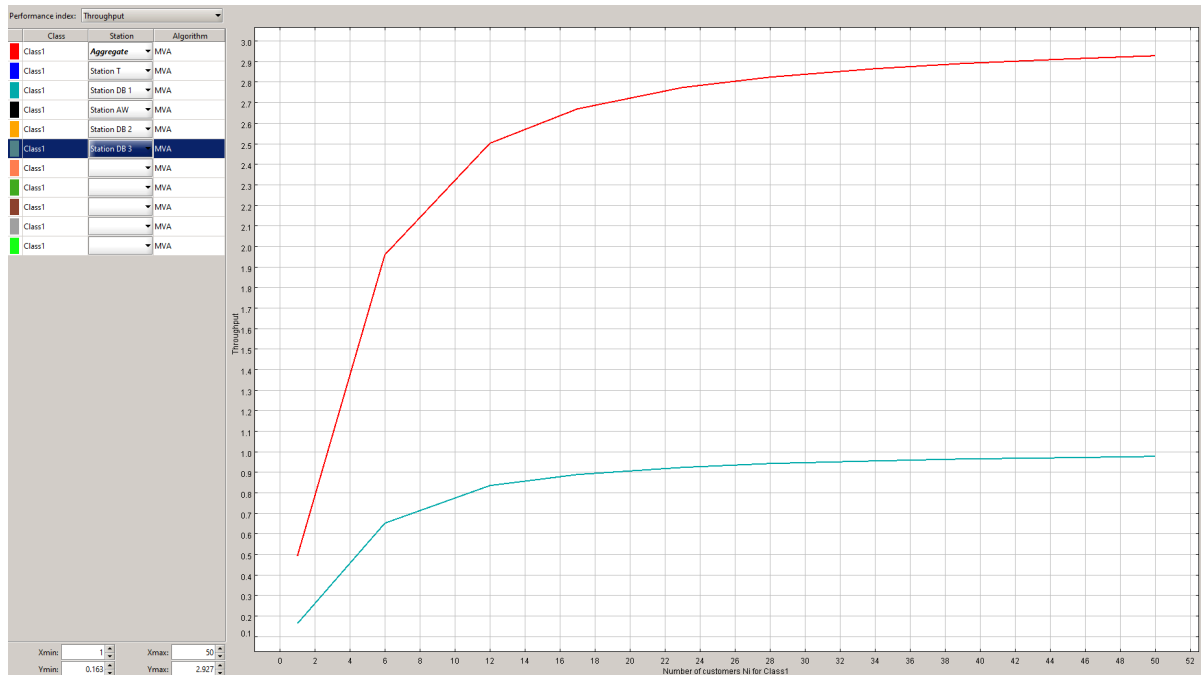
Università
Ca' Foscari
Venezia

Right below the results in graphs:





Università Ca' Foscari Venezia



As the reader can see, the throughput with three database stations increases by nearly a factor of three; similarly, the response time decreases by approximately the same factor.