



Università
Ca'Foscari
Venezia

Scalability and Performance Analysis of a Web Application

Students:

Federico Carraro, 847386

Andrea Brion, 860595

Academic year:

2024/2025



Università
Ca'Foscari
Venezia

Introduction & project's goal

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.

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?

The goal:

We are interested in computing the speed of the system, the throughput in job/seconds as well as the response time, measured in seconds, which is crucial from the point of view of a user.

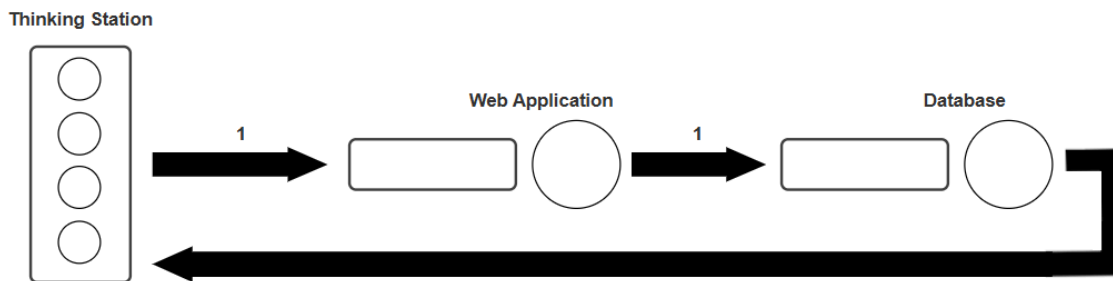
Another important feature to take into account is the system's scalability focusing on its load factor when the number of jobs grows over time and consequently we need to identify the so-called bottleneck responsible to limit its performance.



Università
Ca' Foscari
Venezia

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 (single-class). In addition, each station is always visited by every user so the relative visit ratio for each one is equal to 1.

The throughput with a N amount of users $X(N)$ of such a system is defined by the following formula:

$$X(N) = N / E(T[N])$$

where $E(T[N])$ is the system time



Università
Ca' Foscari
Venezia

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

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

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



Università Ca' Foscari Venezia

Thread Group

Name: Thread Group

Comments:

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 10000

☐ Same user on each iteration

☐ Delay Thread creation until needed

☐ Specify Thread lifetime

Duration (seconds):

Startup delay (seconds):

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.

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
------	-------	-------------	--------------

The JDBC configuration element, where it's possible to see the database's connection settings.
We've used a local database in the configuration, in fact the URL specified is a local one



Università Ca' Foscari Venezia

JDBC Request

Name:

Commento:

Variable Name Bound to Pool:

Variable Name of Pool declared in JDBC Connection Configuration:

SQL Query

Query Type:

Query:

```
1 SELECT DISTINCT
2     nb, 'nb' as nbname,
3     nb, 'nbname' as nbname,
4     tb, 'tb' as tbname,
5     tb, 'tbname' as tbname,
6     tb, 'tbname' as tbname,
7 FROM   nb
8 JOIN LATERAL unnest(string_to_array(nb, ' ')) AS known_title(known) ON true
9 JOIN 'title' ON tb ON tb, known = known_title.known
10 WHERE nb, 'nbname' IN (
11     'Albert Einstein',
12     'Albert Einstein',
13     'Albert Einstein',
14     'Albert Einstein',
15     'Albert Einstein',
16     'Albert Einstein',
17     'Albert Einstein',
18     'Albert Einstein',
19     'Albert Einstein',
20     'Albert Einstein',
21 )
```

Parameter values:

Parameter types:

Variable names:

Result variable name:

Query timeout (s):

Limit ResultSet:

Handle ResultSet:

The JDBC sampler.

Thread Group

Name:

Comments:

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):

Ramp-up period (seconds):

Loop Count: ☒ Infinite

☐ Same user on each iteration

☐ Delay Thread creation until needed

☐ Specify Thread lifetime

Duration (seconds):

Startup delay (seconds):

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.



Università Ca' Foscari Venezia

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.

Flow Control Action

Name: Think Time

Comments:

Logical Action on Thread

☒ Pause Duration (milliseconds):

☐ Start Next Thread Loop ☐ Go to next iteration of Current Loop

☐ Break Current Loop

Logical Action on Thread/Test

☐ Stop ☐ Stop Now

Target: Current Thread

Here we see Flow Control Action.



Università
Ca' Foscari
Venezia

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

jMeter analysis

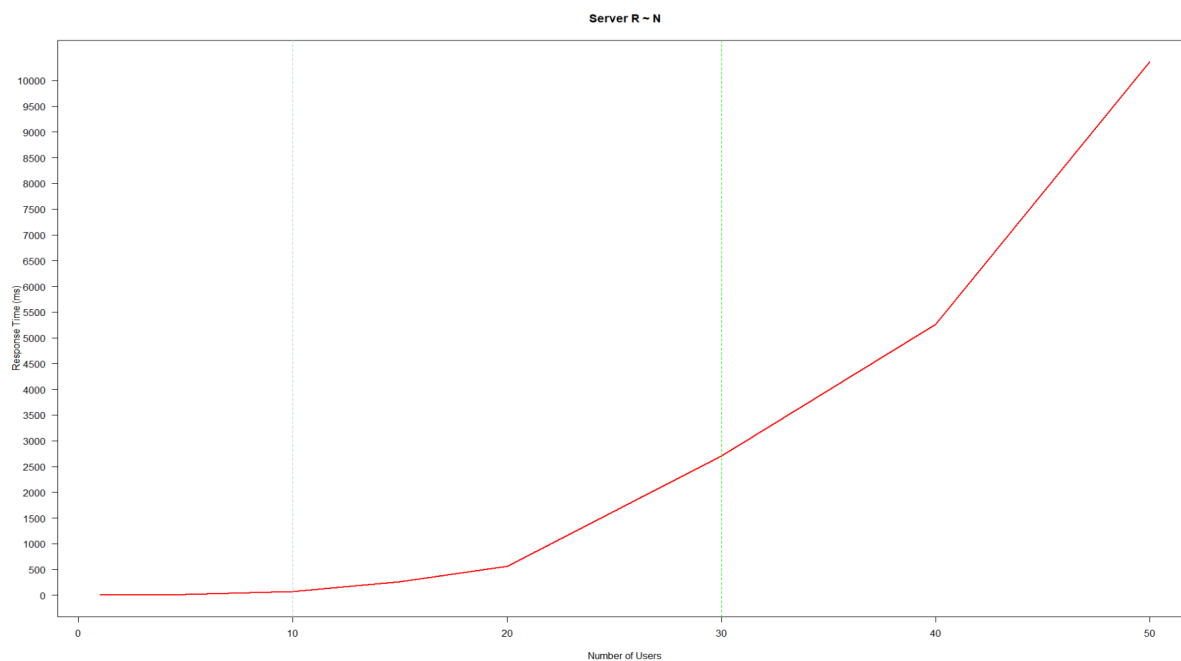
The results of **jMeter** obtained with the increase of the number of customer.
To better readability we've used rStudio to make the plots.
Analyzing the results we've found that

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

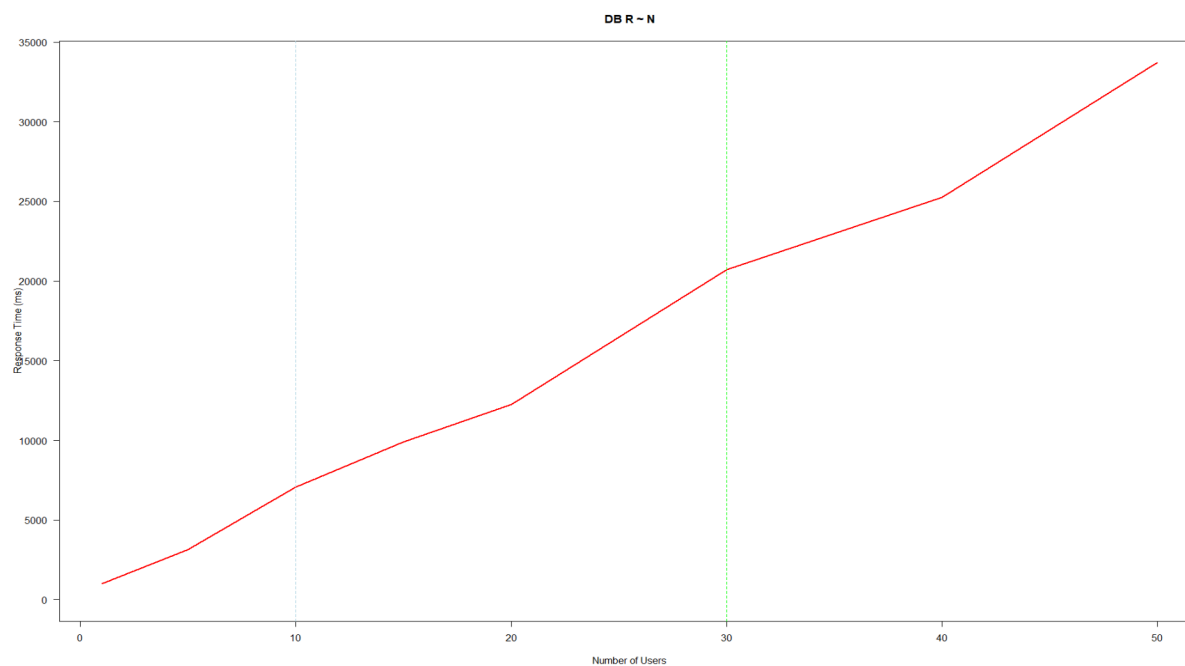
As indicated by the dashed lines in the graphs that follows:



Università
Ca' Foscari
Venezia



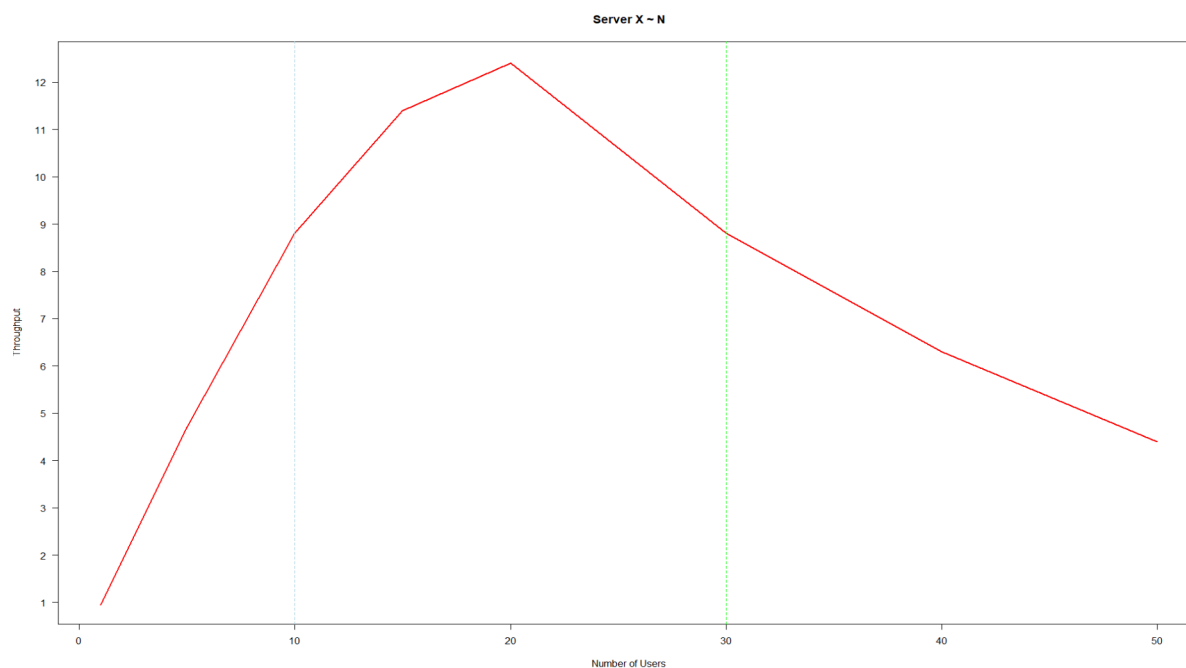
The response time of the server in the system not optimized.



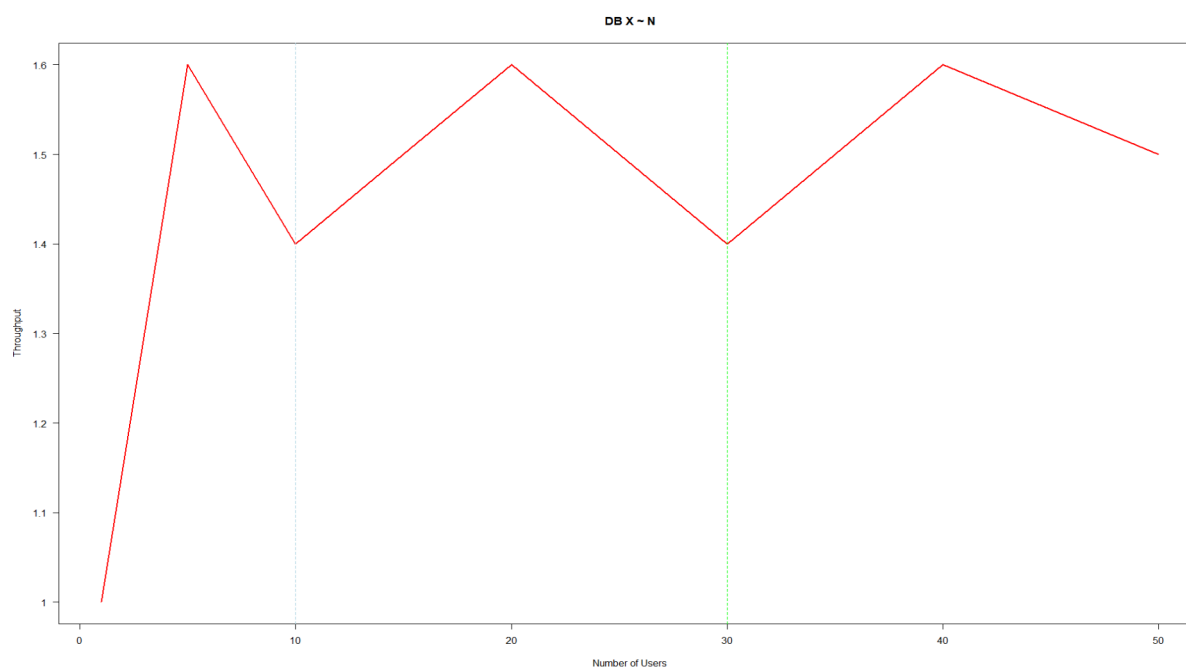
The response time of the database in the system not optimized.



Università
Ca' Foscari
Venezia



The throughput of the server in the system not optimized.



The throughput of the database in the system not optimized.



Università
Ca' Foscari
Venezia

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:

*	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:

*	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.

*	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.

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 (N1): 50

Steps (n. of executions): 10

The What-If.



Università
Ca' Foscari
Venezia

Bottleneck analysis, traffic equations, forced flow law, plots

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: $E[V_i] = 1$, $i = 1, 2, 3$

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!



Università
Ca' Foscari
Venezia

Optimal number of users

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 $\mu_{\text{bottleneck}}/E[V_{\text{bottleneck}}] = 1/0.9830 = 1.0173$ job/s

which imposes the following asymptotic limits to our interactive system:

Lower bound for the expected response time:

when a single job is in the system, we have that:

$$E[R] \geq \max(E[D_{\text{system}}], N \cdot E[D_{\text{bottleneck}}] - E[Z])$$

where $E[D_{\text{system}}] = E[D_{\text{app}}] + E[D_{\text{DB}}]$

Upper bound for the throughput:

while for the throughput:

$$X \leq \min(N / (E[D_{\text{system}}] + E[Z]), 1 / E[D_{\text{bottleneck}}])$$

Using one of the two bounds we can compute the optimal number of users:

$$N_{\text{opt}} = (E[D_{\text{system}}] + E[Z]) / E[D_{\text{DB}}] = (0.009 + 0.9830 + 1.050) / 0.9830 = 2.0773 \sim 2 \text{ users}$$

which is the best tradeoff between system throughput, response time and utilization.

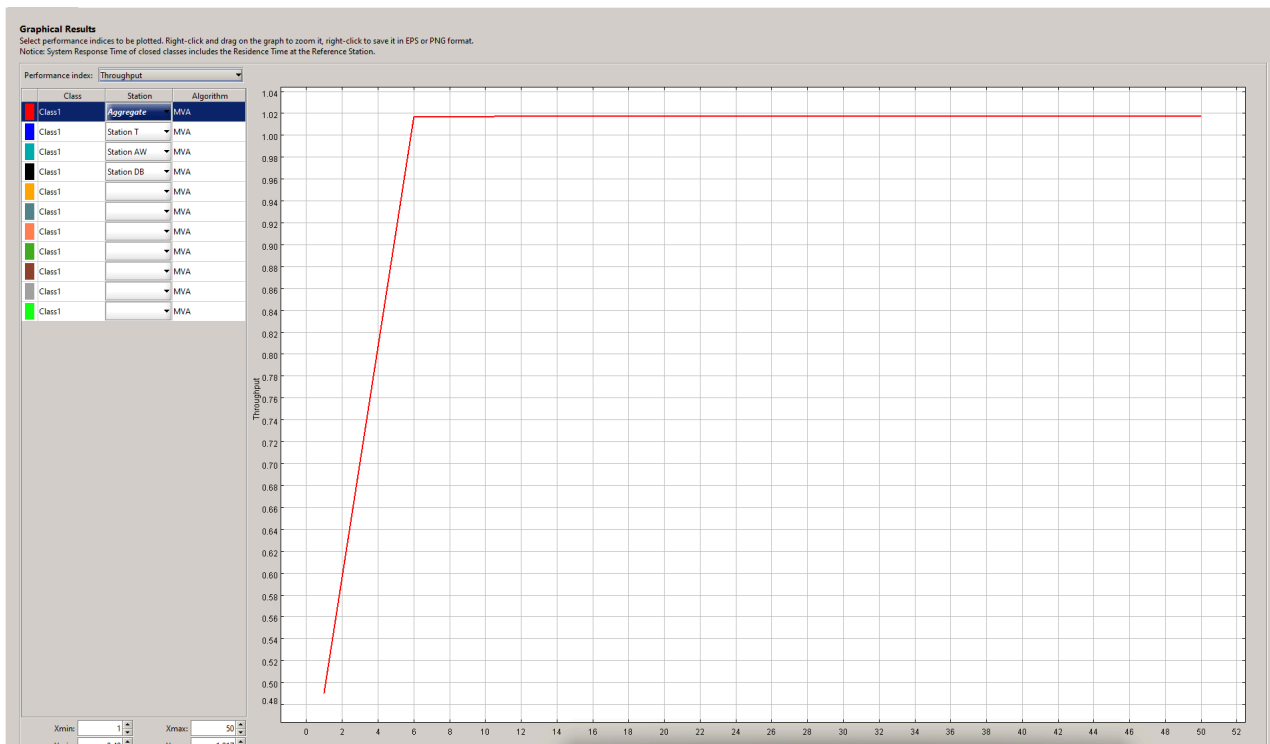
The reader might think that the optimal number of users is too low given the plots below.

However, this number represents the point where the system starts to saturate not where it becomes unstable.

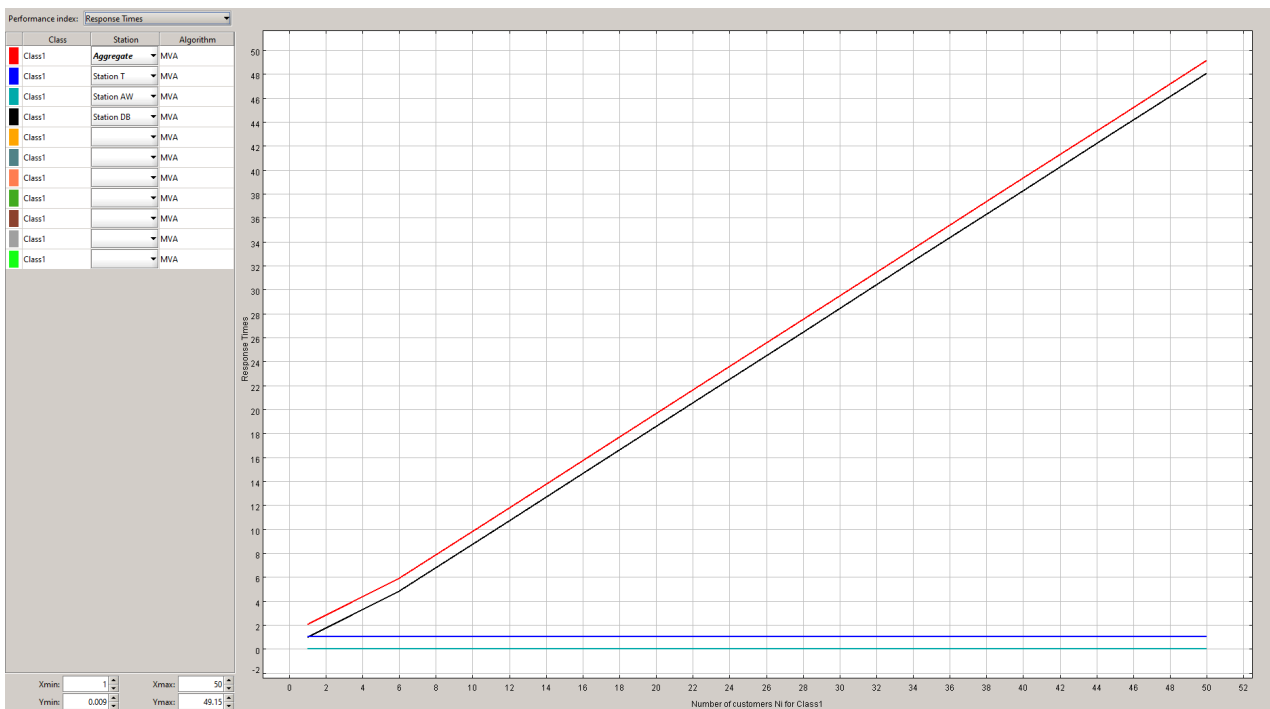
Between 2 and 6 users, the system remains stable but experiences increased latency.



Università
Ca' Foscari
Venezia



Throughput of the system not optimized.



Response time of the system not optimized.



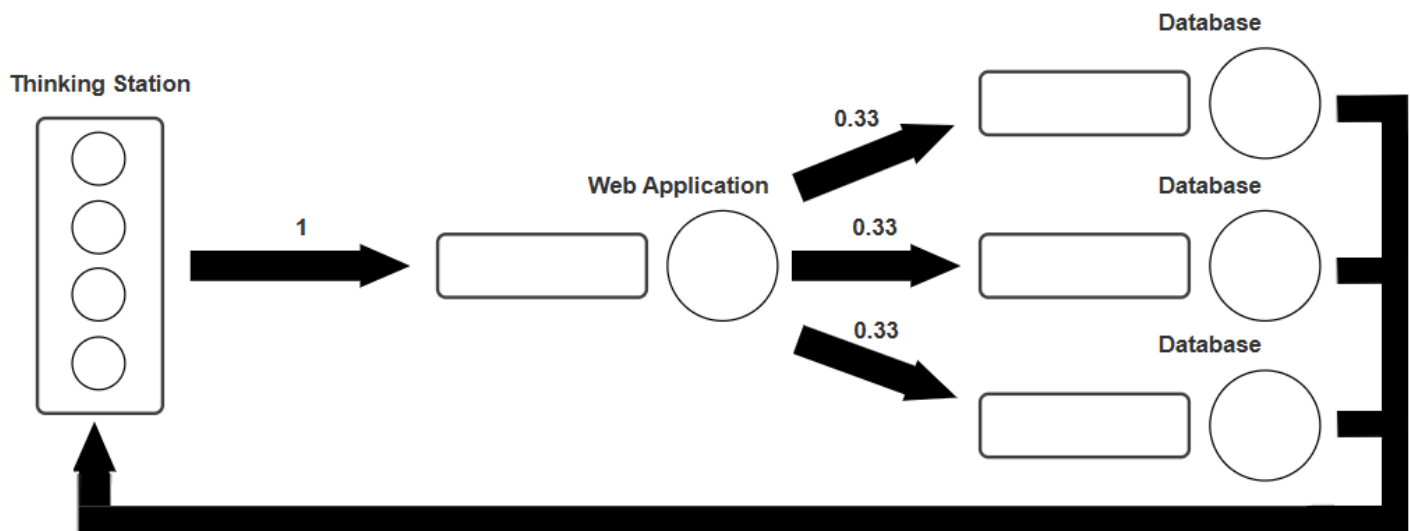
Università
Ca' Foscari
Venezia

System optimization, analysis and results

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.



Let's compute again the bottleneck and the optimal number of users

the web application has relative visit-ratio equal to $E[V_{app}] = 1$ while, now, each database have $E[V_{DBi}] = 1/3$ $i=1,2,3$

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

service demand i-th database: $E[D_{DBi}] = E[V_{DBi}] * 1 / \mu_{DBi} = 0.9830 / 3 = 0.3276$ $i=1,2,3$

still the databases are the bottleneck but now the service demand is lower which speeds up our system so that:

$$X_{system} = \mu_{bottleneck} / E[V_{bottleneck}] = (1/0.9830) * 3 = 1.0173 * 3 = 3.052 \text{ job/s}$$

$$N_{opt} = (E[D_{system}] + E[Z]) / E[D_{DB}] = ((0.009 + 0.3276 + 0.3276 + 0.3276) + (1.050)) / (0.3276) = 6.2332 \sim 6 \text{ users}$$

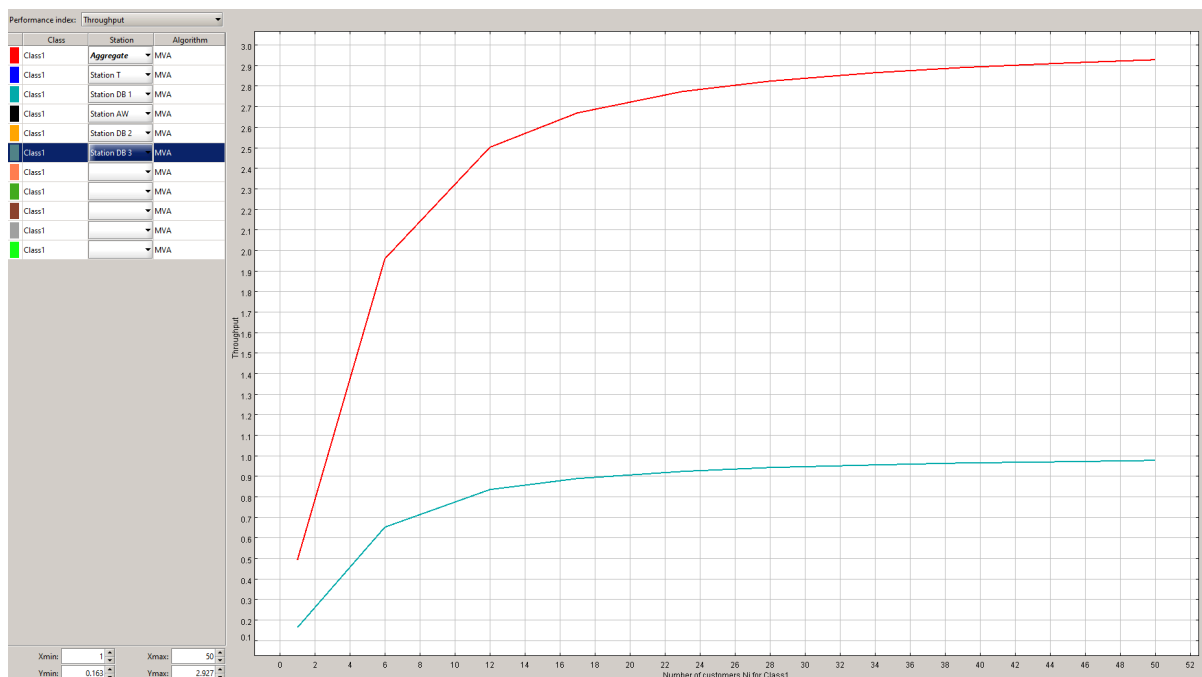
Right below the results in graphs:



Università Ca' Foscari Venezia



Throughput system after optimization.



Throughput system after optimization.

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.