

# Lab 3 – Pygmy, The Book Store

Submission by:

1. Adarsh Kolya - [akolya@umass.edu](mailto:akolya@umass.edu) (Spire ID - 33018261)
2. Vignesh Radhakrishna - [vradhakrishn@umass.edu](mailto:vradhakrishn@umass.edu) (Spire ID - 32577580)
3. Brinda Murulidhara - [bmurulidhara@umass.edu](mailto:bmurulidhara@umass.edu) (Spire ID - 32578418)

## Contents

1. Introduction .....	2
1.1 Goals.....	2
2. Components and Interfaces .....	4
2.1 cache(itemNum) .....	4
2.2 Load balancer component.....	5
2.3 Fault Tolerance .....	5
2.4 Ensuring consistency among the replicas.....	6
2.5 Crash recovery .....	7
3. How to run the system? .....	8
To run the system on localhost .....	8
To run the system on AWS EC2 instances .....	8
Usage .....	9
CLI.....	9
4. Fault tolerance.....	10
Goals of the Experiment .....	10
5. Areas of improvement.....	24
6. Design Trade-offs.....	25
7. Evaluation and Performance .....	26
7.1 Lookups – With and without caching. ....	26
7.2 Buy – With and without caching .....	27

# 1. Introduction

The objective of this project is to implement a multi-tier web application. The web application is a simple version of a book-store wherein users can search for the availability of books by topic or using the unique identifier assigned to every book. Users can then buy books of their choice depending on the availability. REST APIs are exposed to perform basic CRUD operations as part of the backend and frontend microservices. Flask, a light-weight web framework is used to create the distributed web application. SQLite is used as the backend database.

## 1.1 Goals

**1.** Build the backend and frontend components of a two-tier web application wherein microservices are used at each tier. Expose the search, lookup, and buy functionalities as REST APIs in each microservice, wherein the client hits the REST endpoint, and the server services the client request.

**2. Replication** – The order server and catalog servers are replicated, and the load is evenly distributed across these servers. This not only ensures that the time the system needs to process their requests is reduced but also provides a way to continue serving requests even when one or more servers go down.

**3. Caching** – The responses to lookup requests are cached in the frontend server. This reduces the time taken to respond to subsequent requests for the same book item.

**4. Consistency** – The data consistency across replicas is ensured through propagation of requests from primary (server receiving the write request) to replicas. The cache in the front-end server is invalidated every time there is a database update request on an item in the cache.

**5. Fault tolerance** – The system is made fault tolerant using replicas and heartbeat mechanism to identify failed nodes. Client requests are forwarded only to non-faulty nodes and the faulty nodes re-sync data by contacting working replicas once they begin operating normally after a crash.

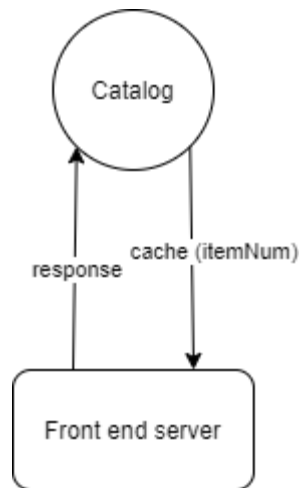
**6. Containerization and docker** – The entire application is bundled and executed in an isolated environment. This ensures that the application with all the required libraries and other dependencies is deployed as one package.

## 2. Components and Interfaces

The front-end server implements the following interfaces in addition to the interfaces provided as part of [Lab 2](#).

### 2.1 `cache(itemNum)`

This is an interface provided to the catalog server to invalidate a specific item in the front-end server's cache when it is updated in the catalog database. We use a server push model to remove stale cache entries. The catalog server invalidates the specific entry in the cache when there is a write request on the database. If the same item is looked up again, the front-end server caches the response. Subsequent lookup requests for the same item do not result in calls to the catalog server as the cached responses are returned. We design this using REST client/server architecture, and a high-level design is as below:



*Figure 1 - Implementation details of the cache invalidation interface*

## 2.2 Load balancer component

The front-end server acts as a load balancer and evenly distributes the incoming requests across the order and catalog servers using round-robin technique.

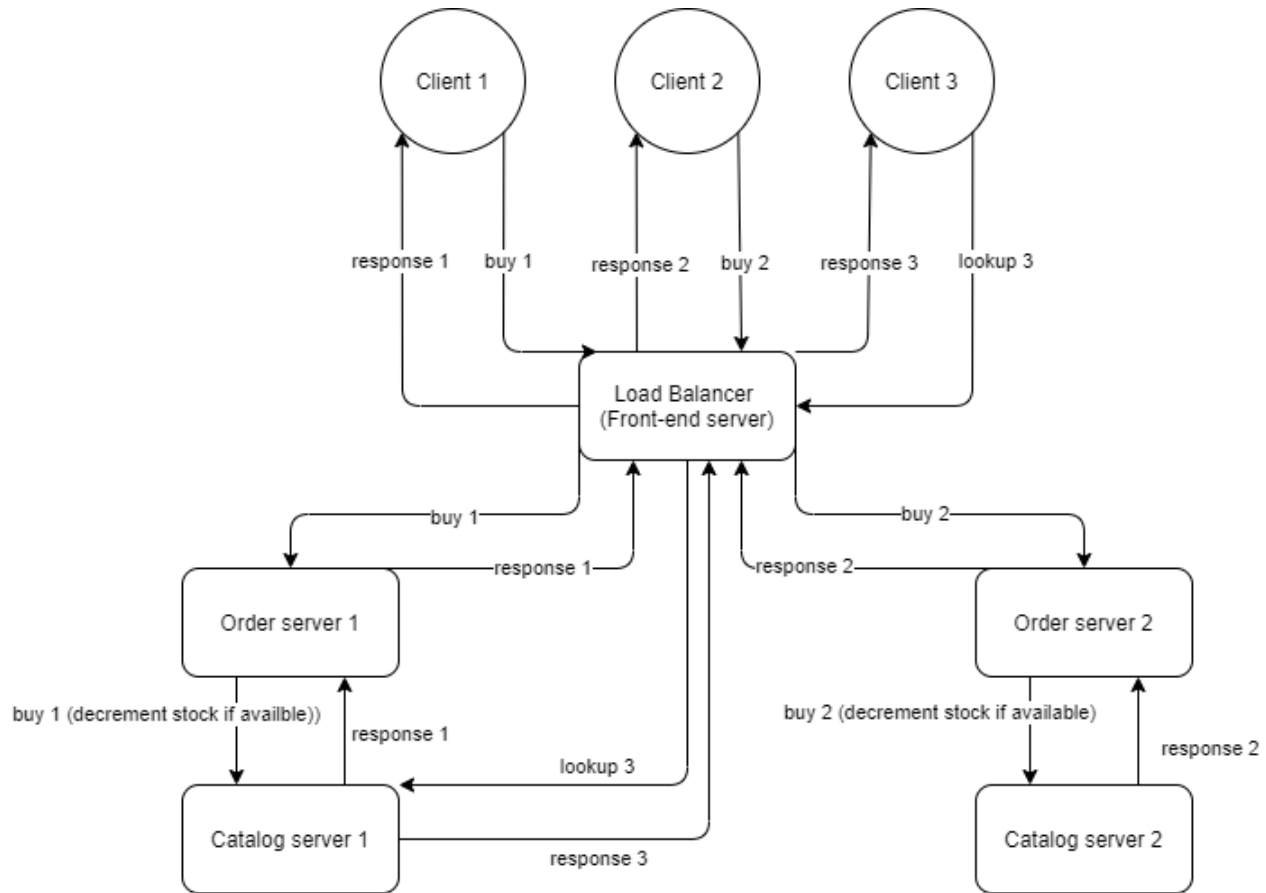


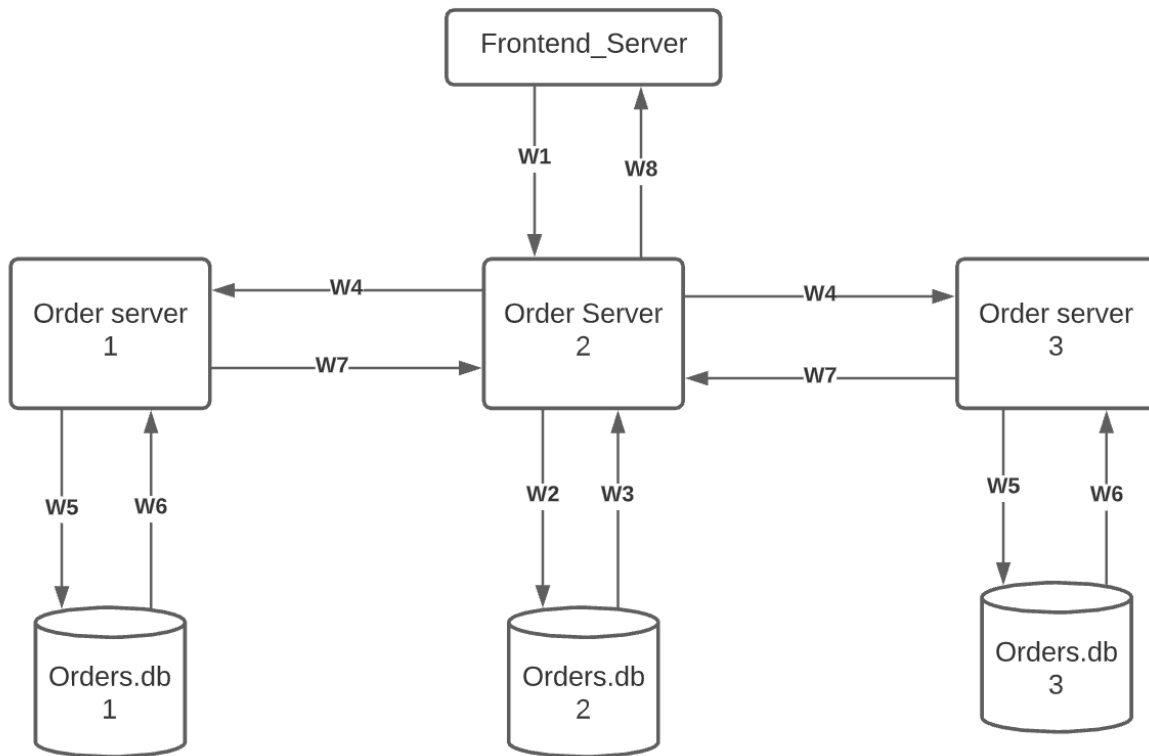
Figure 2 - Load balancer component

## 2.3 Fault Tolerance

The front-end server checks if the replicas are alive by sending heart beats requests at regular intervals. Depending on the status of the replica servers, the front-end server performs load balancing and forwards requests to non-fault replicas.

## 2.4 Ensuring consistency among the replicas

To ensure that the order server and catalog server replicas are in sync with each other, we have followed **Local Write with Replication** protocol among the replicas. In this model the primary copy migrates between the process that want to perform the write. Once the writes have been performed on the primary, this is propagated onto the other replicas as well as shown below:



W1 --> Buy request from the front-end server.  
W2 --> Write to the new primary for orders.db i.e, Orders.db 2.  
W3 --> Response for W2.  
W4 --> Propagation of writes to the replica from the primary server.  
W5 --> Writes on the replica databases.  
W6 --> Response for W5.  
W7 --> Response for W4.  
W8 --> Response for W1.

Figure 3 - Local write with replication protocol to ensure consistency among replicas

## 2.5 Crash recovery

When a replica crashes and then later comes back online, we have a crash recovery mechanism that will help the crashed replica come in sync with the other replicas that were online. The backend servers will have 2 states – INIT and RUNNING. Every server will be in INIT state when it comes up for the first time. As soon as a write is performed on the server's DB, the state is changed to RUNNING. RUNNING state simply means that the server has performed some writes.

Crash recovery logic makes use of these states of the servers and a server synchronizes itself after a crash as described below:

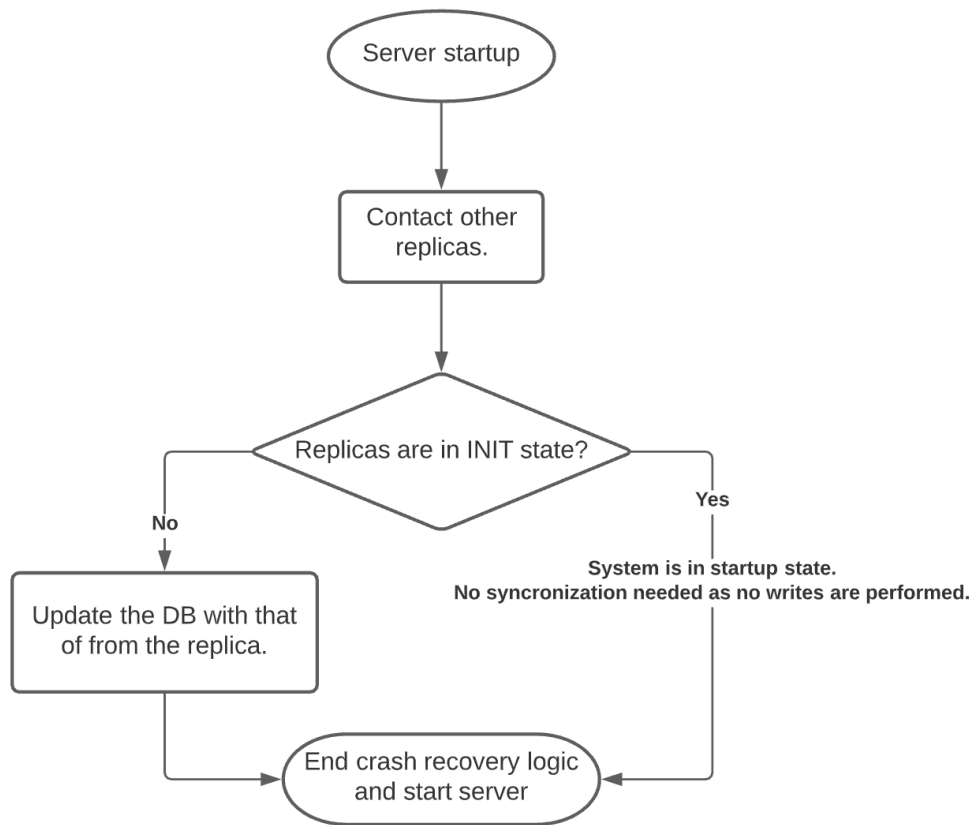


Figure 4 - Crash recovery mechanism

### 3. How to run the system?

#### To run the system on localhost

Run 2 catalog servers, 2 order servers and 1 frontend server on localhost

```
cd Lab-3-Pygmy-The-book-store
bash run.sh
```

#### To run the system on AWS EC2 instances

Create a config file. This file should have 3 lines.

1st line will contain comma separated list of DNS name/IP of servers where catalog server replicas should run.

2nd line will contain comma separated list of DNS name/IP of servers where order server replicas should run.

3rd line will contain DNS name/IP of server where frontend server should run.

*Example :*

```
cat config
ec2-52-204-17-35.compute-1.amazonaws.com,ec2-54-146-69-77.compute-1.amazonaws.com,ec2-54-152-5-125.compute-1.amazonaws.com
ec2-54-157-158-136.compute-1.amazonaws.com,ec2-52-23-179-14.compute-1.amazonaws.com,ec2-107-23-14-212.compute-1.amazonaws.com
ec2-54-237-55-198.compute-1.amazonaws.com
```

With the above config file 3 replicas of catalog server will run on servers listed on line 1. 3 replicas of order server will run on servers listed on line 2. And front end server will run on server on line 3.

Please ensure there is no extra space in between.

Login to a server that has passwordless ssh set up to all the above servers. Clone the repo and do the below :

```
cd Lab-3-Pygmy-The-book-store
bash run.sh config
```



# Usage

## CLI

Open a new shell and execute the below commands.

```
cd LAB-3-PYGMY-THE-BOOK-STORE
# For linux based
pip install virtualenv
virtualenv .venv
source .venv/bin/activate

pip install -r src/requirements.txt
source .venv/bin/activate

cd src/cli
```

### *Lookup*

```
python main.py --frontend_server <ip_of_frontend_server> lookup <item number>
```

Note: --frontend\_server is optional. If not provided, it takes localhost as default.

### *Example*

```
python main.py --frontend_server ec2-35-175-129-185.compute-1.amazonaws.com lookup 1
```

### *Search*

```
python main.py --frontend_server <ip_of_frontend_server> search --topic "<topic>"
```

### *Example*

```
python main.py search --topic "distributed systems"
```

### *Buy*

```
python main.py --frontend_server <ip_of_frontend_server> buy <item number>
```

### *Example*

```
python main.py buy 1
```

Note: For more details refer “How to run the system?” section in the design document for [Lab 2](#).

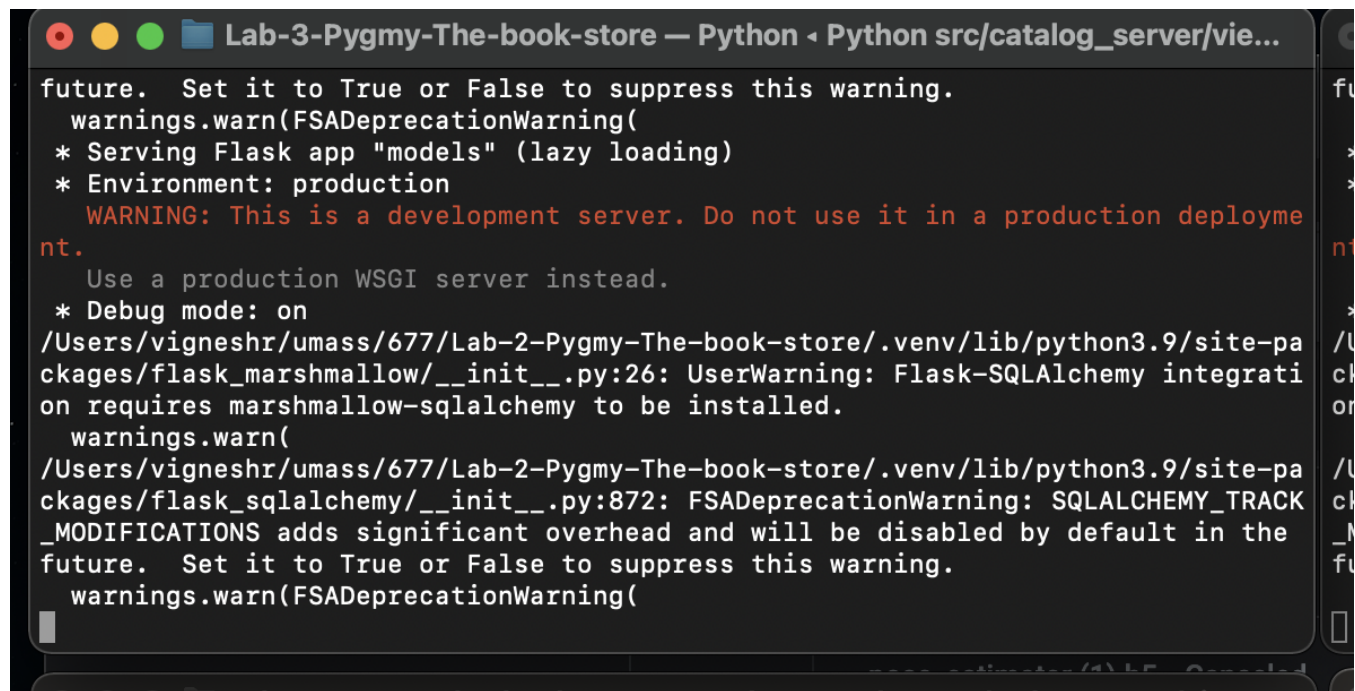
## 4. Fault tolerance

### Goals of the Experiment

1. Show that when one of the catalog server is down and the frontend server sends request to it, the request is forwarded to an active catalog server by the load balancer.
2. Show that when one of the order server is down and the frontend server sends request to it, the request is forwarded to an active order server by the load balancer.
3. Show that when the crashed catalog server recovers, it syncs its data base with the other catalog servers.
4. Show that when the crashed order server recovers, it syncs its data base with the other order servers.
5. Show that the frontend server periodically checks the health of the backend servers.

### Sequence of operations for 1:

- 1.1 Ensure all the servers are running.



```
Lab-3-Pygm-The-book-store — Python Python src/catalog_server/vie...
future. Set it to True or False to suppress this warning.
  warnings.warn(FSAdeprecationWarning(
* Serving Flask app "models" (lazy loading)
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployme
nt.
  Use a production WSGI server instead.
* Debug mode: on
/Users/vigneshr/umass/677/Lab-2-Pygmy-The-book-store/.venv/lib/python3.9/site-pa
ckages/flask_marshmallow/__init__.py:26: UserWarning: Flask-SQLAlchemy integrati
on requires marshmallow-sqlalchemy to be installed.
  warnings.warn(
/Users/vigneshr/umass/677/Lab-2-Pygmy-The-book-store/.venv/lib/python3.9/site-pa
ckages/flask_sqlalchemy/__init__.py:872: FSAdeprecationWarning: SQLALCHEMY_TRACK
_MODIFICATIONS adds significant overhead and will be disabled by default in the
future. Set it to True or False to suppress this warning.
  warnings.warn(FSAdeprecationWarning(
```

- 1.2 Trigger 4 search requests to the frontend server.

The screenshot shows a REST client interface with a GET request to `http://localhost:5002/books?topic=distributed systems...`. The request is sent, and the response is displayed in the 'Body' tab. The response is a JSON array of two book objects.

KEY	VALUE	DESCRIPTION
<input checked="" type="checkbox"/> topic	distributed systems	
Key	Value	Description

```
[
  {
    "id": 1,
    "title": "How to get a good grade in 677 in 20 minutes a day"
  },
  {
    "id": 2,
    "title": "RPCs for Dummies"
  }
]
```

Status: 200 OK Time: 53 ms Size: 254 B Save Response

1.3 Check the frontend.log and note that the requests are evenly distributed among catalog servers (5000 and 5003).

```
frontend.log U X
frontend.log
1 Trying to connect to http://localhost:5000/books
2 Trying to connect to http://localhost:5003/books
3 Trying to connect to http://localhost:5000/books
4 Trying to connect to http://localhost:5003/books
5
```

1.4 Stop the first catalog server (5000).

Lab-3-PygmY-The-book-store — -zsh — 80x18

```
warnings.warn(FSADeprecationWarning(
* Serving Flask app "models" (lazy loading)
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployment.
  Use a production WSGI server instead.
* Debug mode: on
/Users/vigneshr/umass/677/Lab-2-PygmY-The-book-store/.venv/lib/python3.9/site-packages/flask_marshmallow/__init__.py:26: UserWarning: Flask-SQLAlchemy integration requires marshmallow-sqlalchemy to be installed.
  warnings.warn(
/Users/vigneshr/umass/677/Lab-2-PygmY-The-book-store/.venv/lib/python3.9/site-packages/flask_sqlalchemy/__init__.py:872: FSADeprecationWarning: SQLALCHEMY_TRACK_MODIFICATIONS adds significant overhead and will be disabled by default in the future. Set it to True or False to suppress this warning.
  warnings.warn(FSADeprecationWarning(
^C%
(.venv) (base) vigneshr@Vigneshs-MacBook-Pro Lab-3-PygmY-The-book-store %
```

### 1.5 Trigger 4 search requests to the frontend server.

The screenshot shows a REST client interface. At the top, the method is **GET** and the URL is `http://localhost:5002/books?topic=distributed systems...`. Below the URL bar, there are tabs for **Params**, **Authorization**, **Headers (6)**, **Body**, **Pre-request Script**, **Tests**, and **Settings**. The **Params** tab is active, showing a table of query parameters:

	KEY	VALUE	DESCRIPTION
<input checked="" type="checkbox"/>	topic	distributed systems	
	Key	Value	Description

Below the table, there are tabs for **Body**, **Cookies**, **Headers (4)**, and **Test Results**. The **Body** tab is active, showing the response in **JSON** format:

```
1 {
2   {
3     "id": 1,
4     "title": "How to get a good grade in 677 in 20 minutes a day"
5   },
6   {
7     "id": 2,
8     "title": "RPCs for Dummies"
9   }
10 }
```

### 1.6 Check the frontend.log and note an exception log saying the first catalog server is down.

Also note that the subsequent 4 requests are all sent to the other catalog server (5003).

```
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5000/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5000/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5000/books
Exception occurred. HTTPConnectionPool(host='localhost', port=5000): Max retries exceeded with url: /books?topic=distributed+system
The catalog server with ip http://localhost:5000 seems to be down. Will retry the request with the other catalog servers.
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5003/books
```

## Sequence of operations for 2:

### 2.1 Ensure all the servers are running.

```
Lab-3-Pygmy-The-book-store — Python ▶ Python src/catalog_server/vie...

future. Set it to True or False to suppress this warning.
warnings.warn(FSAdeprecationWarning(
* Serving Flask app "models" (lazy loading)
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployment.
  Use a production WSGI server instead.
* Debug mode: on
/Users/vigneshr/umass/677/Lab-2-Pygmy-The-book-store/.venv/lib/python3.9/site-packages/flask_marshmallow/__init__.py:26: UserWarning: Flask-SQLAlchemy integration requires marshmallow-sqlalchemy to be installed.
warnings.warn(
/Users/vigneshr/umass/677/Lab-2-Pygmy-The-book-store/.venv/lib/python3.9/site-packages/flask_sqlalchemy/__init__.py:872: FSAdeprecationWarning: SQLALCHEMY_TRACK_MODIFICATIONS adds significant overhead and will be disabled by default in the future. Set it to True or False to suppress this warning.
warnings.warn(FSAdeprecationWarning(
```

### 2.2 Trigger 4 buy requests to the frontend server.

Frontend server buy

POST http://localhost:5002/books/1 Send

Params Authorization Headers (7) Body Pre-request Script Tests Settings Cookies </>

Query Params

KEY	VALUE	DESCRIPTION	...	Bulk Edit
Key	Value	Description		

Body Cookies Headers (4) Test Results Status: 200 OK Time: 43 ms Size: 222 B Save Response

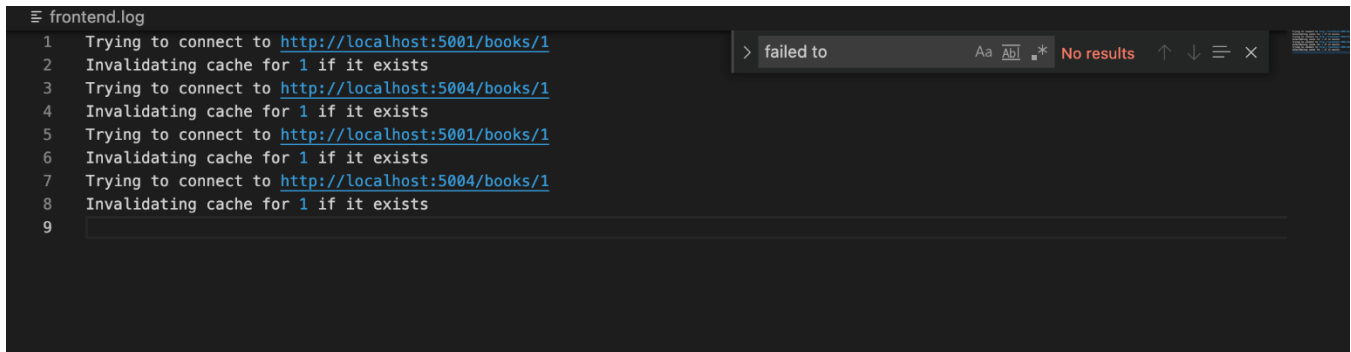
Pretty Raw Preview Visualize JSON

```
1  {"Date": "Apr-30-2021 17:44:42",
2  "ID": 1,
3  "Order status": "Success",
4  "OrderId": 79
5  }
```

**200 OK**

Standard response for successful HTTP requests. The actual response will depend on the request method used. In a GET request, the response will contain an entity corresponding to the requested resource. In a POST request the response will contain an entity describing or containing the result of the

2.3 Check the frontend.log and note that the requests are evenly distributed among order servers (5001 and 5004).

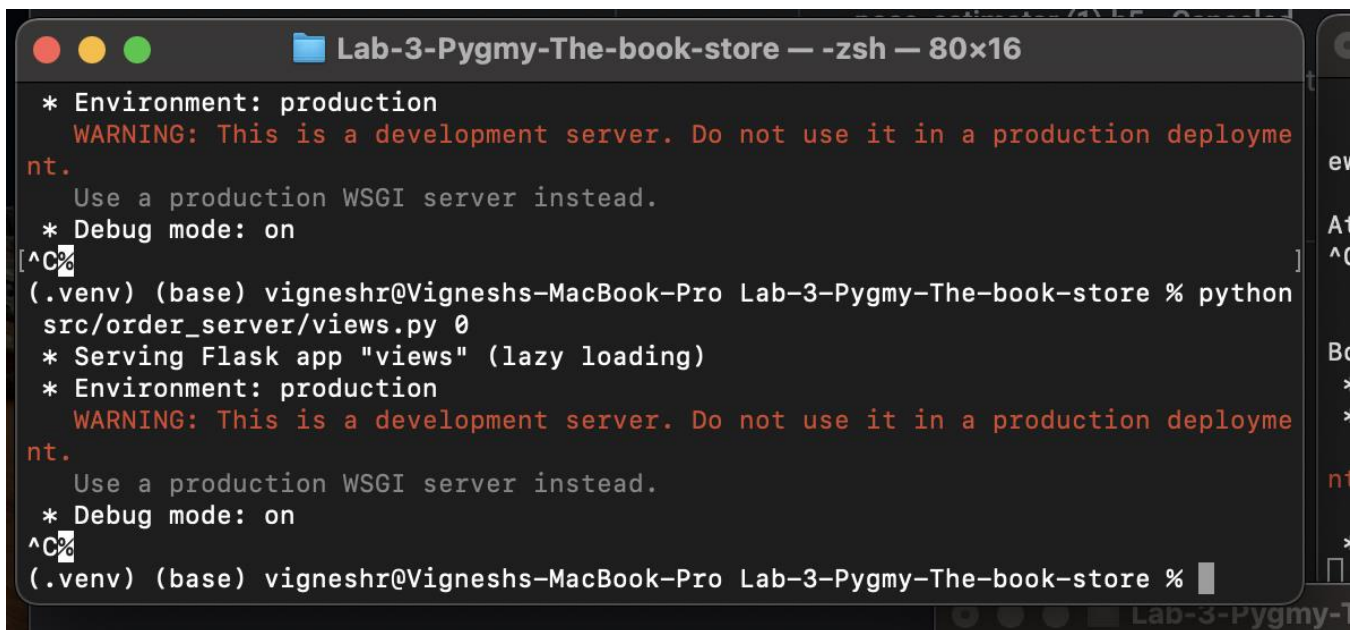


The screenshot shows a log viewer window titled 'frontend.log'. The log contains the following entries:

```
1 Trying to connect to http://localhost:5001/books/1
2 Invalidating cache for 1 if it exists
3 Trying to connect to http://localhost:5004/books/1
4 Invalidating cache for 1 if it exists
5 Trying to connect to http://localhost:5001/books/1
6 Invalidating cache for 1 if it exists
7 Trying to connect to http://localhost:5004/books/1
8 Invalidating cache for 1 if it exists
9
```

A search bar at the top right shows the text '> failed to' and 'No results'.

2.4 Stop the first order server (5001)



The screenshot shows a terminal window titled 'Lab-3-Pygm-The-book-store — zsh — 80x16'. The terminal output is as follows:

```
* Environment: production
WARNING: This is a development server. Do not use it in a production deployment.
Use a production WSGI server instead.
* Debug mode: on
[ ^C%
(.venv) (base) vigneshr@Vigneshs-MacBook-Pro Lab-3-Pygm-The-book-store % python
src/order_server/views.py 0
* Serving Flask app "views" (lazy loading)
* Environment: production
WARNING: This is a development server. Do not use it in a production deployment.
Use a production WSGI server instead.
* Debug mode: on
[ ^C%
(.venv) (base) vigneshr@Vigneshs-MacBook-Pro Lab-3-Pygm-The-book-store %
```



## 2.5 Trigger 4 buy requests to the frontend server.

The screenshot shows a REST client interface with the following details:

- Method:** POST
- URL:** http://localhost:5002/books/1
- Buttons:** Send, Cookies, </>, ⓘ
- Query Params:** A table with columns KEY, VALUE, and DESCRIPTION. It contains one row with 'Key', 'Value', and 'Description'.
- Body:** A tabbed interface with 'Body', 'Cookies', 'Headers (4)', and 'Test Results'. The 'Body' tab is active, showing a JSON response in 'Pretty' format:

```
1 {
2   "Date": "Apr-30-2021 17:44:42",
3   "ID": 1,
4   "Order status": "Success",
5   "OrderId": 79
6 }
```
- Status Bar:** Status: 200 OK, Time: 43 ms, Size: 222 B, Save Response
- 200 OK Popover:** A tooltip explaining the status code: "Standard response for successful HTTP requests. The actual response will depend on the request method used. In a GET request, the response will contain an entity corresponding to the requested resource. In a POST request the response will contain an entity describing or containing the result of the"

## 2.6 Check the frontend.log and note an exception log saying the first order server is down.

Also note that the subsequent 4 requests are all sent to the other order server (5003).



```
1 Trying to connect to http://localhost:5001/books/1
2 Invalidating cache for 1 if it exists
3 Trying to connect to http://localhost:5004/books/1
4 Invalidating cache for 1 if it exists
5 Trying to connect to http://localhost:5001/books/1
6 Exception occurred. HTTPConnectionPool(host='localhost', port=5001): Max retries exce
7 The order server with ip http://localhost:5001 seems to be down. Will retry the requ
8 Trying to connect to http://localhost:5004/books/1
9 Invalidating cache for 1 if it exists
10 Trying to connect to http://localhost:5004/books/1
11 Invalidating cache for 1 if it exists
12 Trying to connect to http://localhost:5004/books/1
13 Invalidating cache for 1 if it exists
14 Trying to connect to http://localhost:5004/books/1
15 Invalidating cache for 1 if it exists
16
```

### Sequence of operations for 3:

3.1 Ensure all the servers are running.

```
Lab-3-Pygm-The-book-store — Python Python src/catalog_server/vie...
future. Set it to True or False to suppress this warning.
  warnings.warn(FSAdeprecationWarning(
* Serving Flask app "models" (lazy loading)
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployme
nt.
  Use a production WSGI server instead.
* Debug mode: on
/Users/vigneshr/umass/677/Lab-2-Pygm-The-book-store/.venv/lib/python3.9/site-pa
ckages/flask_marshmallow/__init__.py:26: UserWarning: Flask-SQLAlchemy integrati
on requires marshmallow-sqlalchemy to be installed.
  warnings.warn(
/Users/vigneshr/umass/677/Lab-2-Pygm-The-book-store/.venv/lib/python3.9/site-pa
ckages/flask_sqlalchemy/__init__.py:872: FSAdeprecationWarning: SQLALCHEMY_TRACK
_MODIFICATIONS adds significant overhead and will be disabled by default in the
future. Set it to True or False to suppress this warning.
  warnings.warn(FSAdeprecationWarning(
```

3.2. Stop the first catalog server (5000)

```
Lab-3-Pygm-The-book-store — -zsh — 80x18
  warnings.warn(FSAdeprecationWarning(
* Serving Flask app "models" (lazy loading)
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployme
nt.
  Use a production WSGI server instead.
* Debug mode: on
/Users/vigneshr/umass/677/Lab-2-Pygm-The-book-store/.venv/lib/python3.9/site-pa
ckages/flask_marshmallow/__init__.py:26: UserWarning: Flask-SQLAlchemy integrati
on requires marshmallow-sqlalchemy to be installed.
  warnings.warn(
/Users/vigneshr/umass/677/Lab-2-Pygm-The-book-store/.venv/lib/python3.9/site-pa
ckages/flask_sqlalchemy/__init__.py:872: FSAdeprecationWarning: SQLALCHEMY_TRACK
_MODIFICATIONS adds significant overhead and will be disabled by default in the
future. Set it to True or False to suppress this warning.
  warnings.warn(FSAdeprecationWarning(
^C%
(.venv) (base) vigneshr@Vigneshs-MacBook-Pro Lab-3-Pygm-The-book-store %
```

3.3 Trigger a lookup request to the frontend server and note the count of the book (1154).

bookstore / Front end server lookup

GET http://localhost:5002/books/1... Send

Params Authorization Headers (6) Body Pre-request Script Tests Settings Cookies

Query Params

KEY	VALUE	DESCRIPTION	...	Bulk Edit
Key	Value	Description		

Body Cookies Headers (4) Test Results Status: 200 OK Time: 19 ms Size: 172 B Save Response

Pretty Raw Preview Visualize JSON

```
1 {
2   "cost": 30.0,
3   "count": 1154
4 }
```

3.4 Note in the frontend.log file that the request was sent to the other catalog server (5003).

```
Trying to connect to http://localhost:5000/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5000/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5000/books
Exception occurred. HTTPConnectionPool(host='localhost', port=5000): Max retr
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5003/books/1
Caching lookup results of 1
```

3.5 Trigger a buy request to the frontend server. The expected count is now 1153.

bookstore / front-end server buy

POST http://localhost:5002/books/1 Send

Params Authorization Headers (7) Body Pre-request Script Tests Settings Cookies

Query Params

KEY	VALUE	DESCRIPTION	Bulk Edit
Key	Value	Description	

Body Cookies Headers (4) Test Results Status: 200 OK Time: 43 ms Size: 222 B Save Response

Pretty Raw Preview Visualize JSON

```

1  {
2    "Date": "Apr-30-2021 17:44:42",
3    "ID": 1,
4    "Order status": "Success",
5    "OrderId": 79
6  }

```

**200 OK**

Standard response for successful HTTP requests. The actual response will depend on the request method used. In a GET request, the response will contain an entity corresponding to the requested resource. In a POST request the response will contain an entity describing or containing the result of the

3.6 Note that the request is again sent to the only catalog server (5003).

```

Trying to connect to http://localhost:5000/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5000/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5000/books
Exception occurred. HTTPConnectionPool(host='localhost', port=5000): Max retr
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5003/books
Trying to connect to http://localhost:5003/books/1
Caching lookup results of 1

```

3.7 Start the first catalog server.

3.8 Check the catalog.log to verify that the catalog server syncs itself by talking to the other catalog server (5003)

```
INFO:werkzeug: * Restarting with stat
INFO:root:Trying to sync with replica with node num 0 - http://localhost:5003
DEBUG:urllib3.connectionpool:Starting new HTTP connection (1): localhost:5003
INFO:werkzeug:127.0.0.1 - - [30/Apr/2021 15:17:51] " [37mGET /table HTTP/1.1 [0m" 200 -
DEBUG:urllib3.connectionpool:http://localhost:5003 "GET /table HTTP/1.1" 200 1036
INFO:root:Replica 0 - http://localhost:5003 is in RUNNING state. We will sync our DB with that of the replica
```

3.9 Trigger a lookup request to the frontend server.

bookstore / Front end server lookup

GET http://localhost:5002/books/1... Send

Params Authorization Headers (6) Body Pre-request Script Tests Settings Cookies </>

Query Params

KEY	VALUE	DESCRIPTION	...	Bulk Edit
Key	Value	Description		

Body Cookies Headers (4) Test Results Status: 200 OK Time: 38 ms Size: 172 B Save Response

Pretty Raw Preview Visualize JSON

```
1 {
2   "cost": 30.0,
3   "count": 1153
4 }
```

3.10 frontend.log shows that the request was now sent to the recovered server since it is back.

3.11 Check the output of the request and verify that the count is as expected (1153) in 3.4.

bookstore / Front end server lookup

GET http://localhost:5002/books/1 ... Send

Params Authorization Headers (6) Body Pre-request Script Tests Settings Cookies </>

Query Params

KEY	VALUE	DESCRIPTION	...	Bulk Edit
Key	Value	Description		

Body Cookies Headers (4) Test Results Status: 200 OK Time: 38 ms Size: 172 B Save Response

Pretty Raw Preview Visualize JSON

```
1 {
2   "cost": 39.0,
3   "count": 1153
4 }
```

#### Sequence of operations for 4:

- 4.1 Ensure all the servers are running.
- 4.2. Stop the first order server (5001).

```
Lab-3-Pygm-The-book-store — -zsh — 80x16
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployment.
  Use a production WSGI server instead.
* Debug mode: on
(.venv) (base) vigneshr@Vigneshs-MacBook-Pro Lab-3-Pygm-The-book-store % python src/order_server/views.py 0
* Serving Flask app "views" (lazy loading)
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployment.
  Use a production WSGI server instead.
* Debug mode: on
(.venv) (base) vigneshr@Vigneshs-MacBook-Pro Lab-3-Pygm-The-book-store %
```



4.3 Trigger a buy request to the frontend server.

4.4 Start the first order server (5001).

4.5 Check the orders.log to verify that the order server (5001) syncs itself by talking to the other order server (5004)

```
INFO:werkzeug: * Running on http://0.0.0.0:5001/ (Press CTRL+C to quit)
INFO:werkzeug: * Restarting with stat
INFO:root:Trying to sync with replica with node num 0 - http://localhost:5004
DEBUG:urllib3.connectionpool:Starting new HTTP connection (1): localhost:5004
INFO:werkzeug:127.0.0.1 - - [30/Apr/2021 17:45:08] " [37mGET /table HTTP/1.1 [0m" 200 -
DEBUG:urllib3.connectionpool:http://localhost:5004 "GET /table HTTP/1.1" 200 5523
INFO:root:Replica 0 - http://localhost:5004 is in RUNNING state. We will sync our DB with that of the replica
```

## Sequence of operations for 5:

5.1 Check the frontend.log for /health endpoint hits and the 200 return code

```
DEBUG:apscheduler.scheduler:Looking for jobs to run
INFO:apscheduler.executors.default:Running job "check_servers (trigger: interval[0:00:30], next run at: 2021-04-30 15:37:07 EDT)"
DEBUG:apscheduler.scheduler:Next wakeup is due at 2021-04-30 15:37:37.384528-04:00 (in 29.994909 seconds)
DEBUG:urllib3.connectionpool:Starting new HTTP connection (1): localhost:5000
DEBUG:urllib3.connectionpool:http://localhost:5000 "GET /health HTTP/1.1" 200 22
DEBUG:urllib3.connectionpool:Starting new HTTP connection (1): localhost:5003
DEBUG:urllib3.connectionpool:http://localhost:5003 "GET /health HTTP/1.1" 200 22
DEBUG:urllib3.connectionpool:Starting new HTTP connection (1): localhost:5001
DEBUG:urllib3.connectionpool:http://localhost:5001 "GET /health HTTP/1.1" 200 22
DEBUG:urllib3.connectionpool:Starting new HTTP connection (1): localhost:5004
DEBUG:urllib3.connectionpool:http://localhost:5004 "GET /health HTTP/1.1" 200 22
INFO:apscheduler.executors.default:Job "check_servers (trigger: interval[0:00:30], next run at: 2021-04-30 15:37:37 EDT)" executed
```

## 5. Areas of improvement

- 1) **Consistency:** We could have better synchronization methods to ensure that the data across all replicas is consistent at all times. The current implementation only ensures eventual consistency.
- 2) **Backup and recovery:** Backups of the databases must be taken regularly to avoid syncing entire database with a non-faulty replica when there is a system failure. The current implementation requires that the entire database is recreated by syncing with a functional replica.



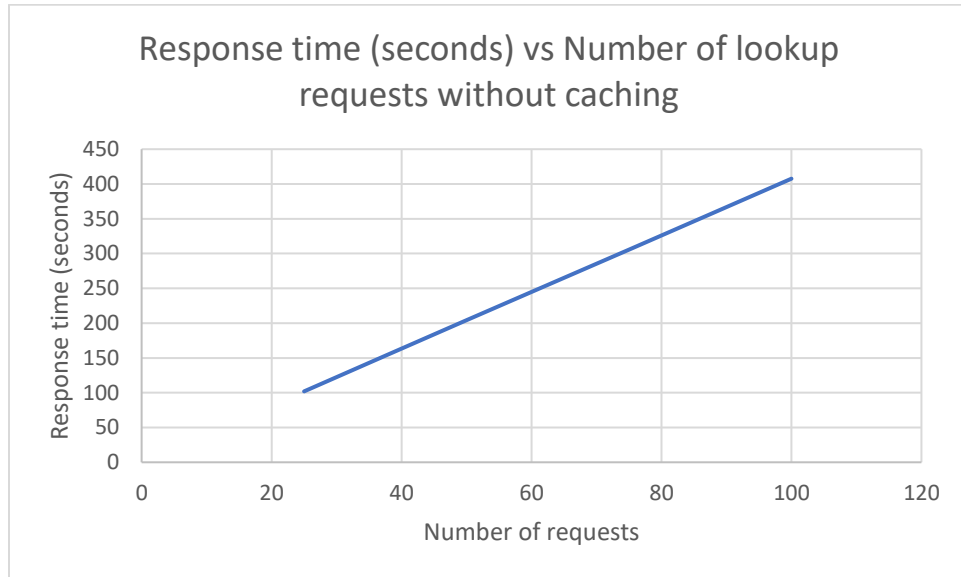
## 6. Design Trade-offs

1) We use Flask as it is lightweight. But Flask is not well known for handling concurrent requests and so as number of clients would increase, the system will become extremely slow resulting in a bad user experience.

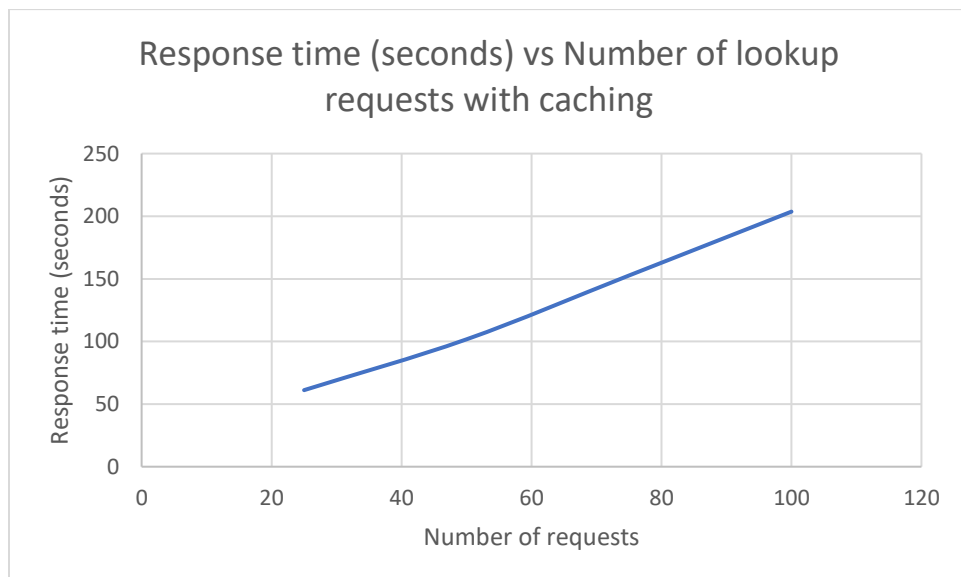
# 7. Evaluation and Performance

## 7.1 Lookups – With and without caching.

Below is a plot showing number of lookup requests on the x axis and time taken to get a response (in seconds) on the y axis. The front-end server does not have any caching mechanism in this case. The average response time is about 4 seconds/request.



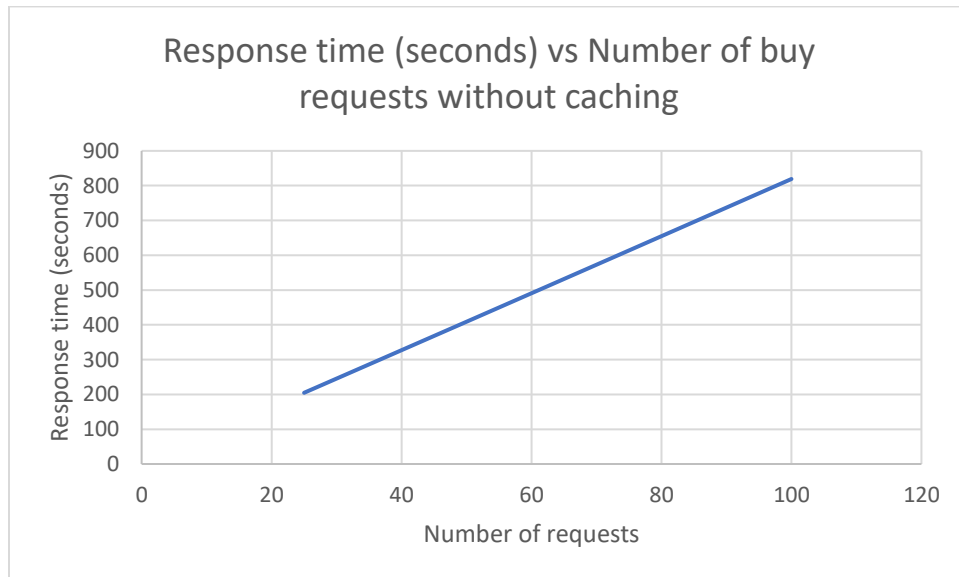
Below is a plot showing number of lookup requests on the x axis and time taken to get a response (in seconds) on the y axis. The front-end server has in-memory caching in this case. The average response time is about 2 seconds/request.



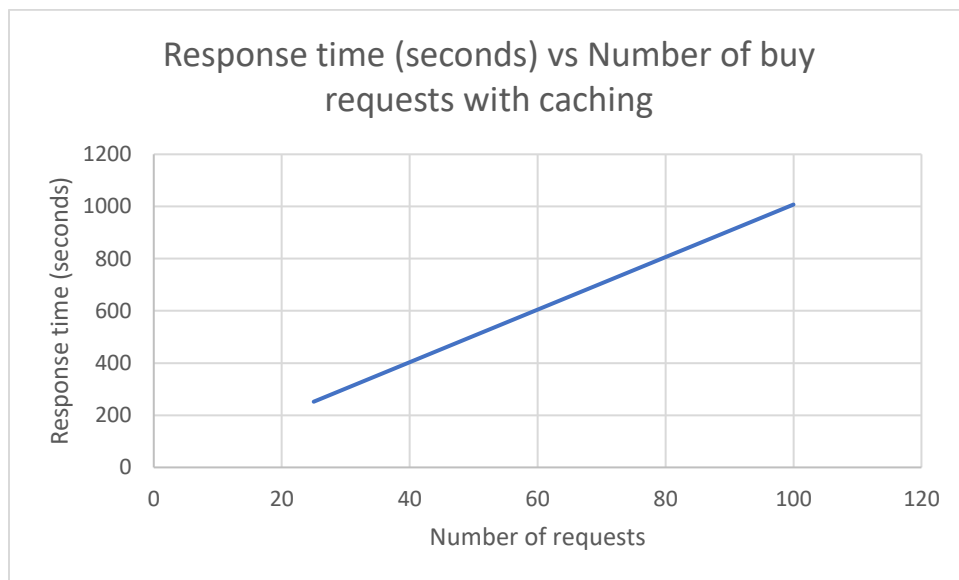
It can be clearly seen from the above two plots that caching significantly reduces response times (nearly reduced by one half).

## 7.2 Buy – With and without caching

Below is a plot showing number of buy requests on the x axis and time taken to get a response (in seconds) on the y axis. The front-end server does not have any caching mechanism in this case. The average response time is about 8.1 seconds/request.

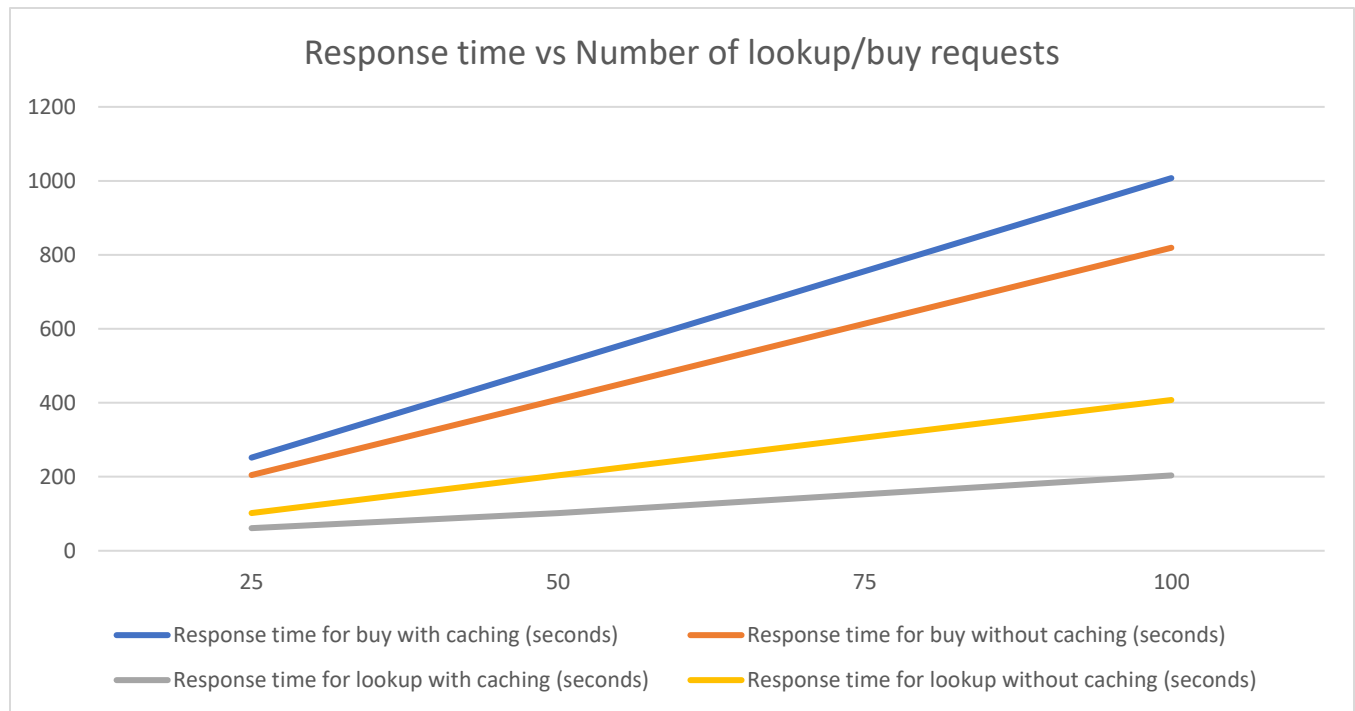


Below is a plot showing number of buy requests on the x axis and time taken to get a response (in seconds) on the y axis. The front-end server has in-memory caching in this case. The average response time is about 10.3 seconds/request.



It can be seen that the cache consistency operations increases the response times slightly (by 2.2 seconds/request on an average) as the buy requests results in an update request on the catalog server which in turn invalidates the cache in the front-end server. However, this latency due to cache operations is not very significant. This maybe since the cache resides in memory and hence the cache invalidate operation does not have any IO which makes is quite fast.

Below is a plot showing trends in both lookup and buy response times with and without caching for comparison.



From the above experiments, we see that the overhead of cache consistency operations is about 1.8 seconds/request.

The latency of a subsequent request if it sees a cache miss is about 4 seconds.