

**An-Najah National University**  
**Faculty of Engineering & Information Technology**

# **Dos-Project - Part 2**

**PhD. Samer Arandi**

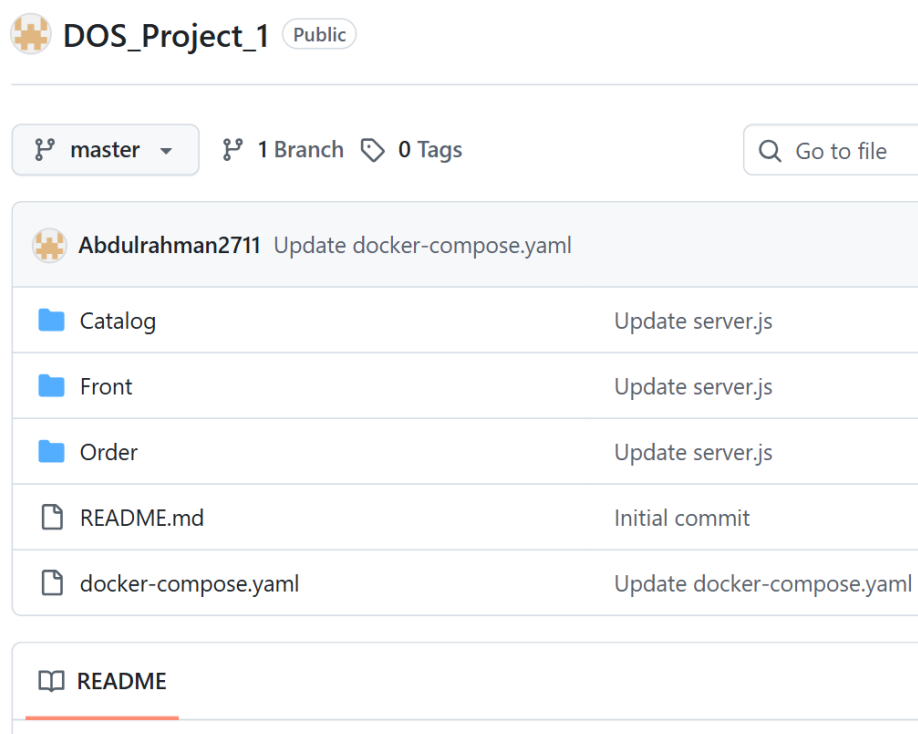
**Abdulrahman Hassoun**  
– 12027846

# 1. Introduction

This report presents the main features implemented in the second phase of the project, which include an in-memory caching system, a load balancing strategy, and the replication of catalog and order services using a microservices architecture with RESTful APIs.

## 2. Project Structure

We have replicated the Catalog service and its database file, and the same was done for the Order service as shown in the following diagram:



## 3. Cache

For read operations, a cache based on a Map structure was implemented to store responses from catalog queries. Cache entries are indexed by the query (for example, `search:<topic>` and `info:<id>`) and contain an expiration timestamp (TTL). When a search request arrives, the frontend first invokes `cacheGet()` to verify if a valid (non-expired) entry exists. In case of a cache hit, the stored response is returned immediately. Otherwise, the request is forwarded to a catalog replica (chosen via round-robin load balancing), the received data is saved with `cacheSet()`, and finally sent back to the client. The cache has a maximum capacity; when full, the oldest entry is evicted. Additionally, a server-push invalidation endpoint (`/cache/invalidate/:id`) is provided so that the catalog service can invalidate cached information about an item before performing write operations, thereby preserving consistency.

```
1 const cache = new Map(); // key -> { value, expiresAt }
2
3 // search => uses the cache and it could hit or mis
```

```

4 app.get("/search/:topic", async (req, res) => [
5   const topic = String(req.params.topic || "");
6   const key = `search:${topic.toLowerCase()}`;
7
8   const hit = cacheGet(key);
9   if (hit) return res.json({ cached: true, data: hit });
10
11   const catalogURL = pickCatalog();
12   const r = await httpJSON(`${catalogURL}/search/${encodeURIComponent
13     (topic)}`);
14   if (!r.ok) return res.status(r.status).json(r.data ?? { error: "
    Catalog error" });
15   cacheSet(key, { property: cached: boolean
16     res.json({ cached: false, replica: catalogURL, data: r.data });
    });

```

## 4. Load Balancing

Requests are distributed across multiple replicas of the backend services. The Round Robin algorithm was chosen for its simplicity and low overhead. For every incoming search request, the frontend picks one replica in a cyclic order, guaranteeing that consecutive requests are sent to different replicas. This method balances the load among the replicas and avoids any single replica becoming a bottleneck.

```

1 const CATALOG_REPLICAS = (process.env.CATALOG_REPLICAS || "http://
   localhost:4001,http://localhost:4003")
2   .split(",")
3   .map(s => s.trim())
4   .filter(Boolean);
5
6 const ORDER_REPLICAS = (process.env.ORDER_REPLICAS || "http://localhost
   :4002,http://localhost:4004")
7   .split(",")
8   .map(s => s.trim())
9   .filter(Boolean);

```

The code that actually implements the load balancing:

```
1 let catIdx = 0;
2 let ordIdx = 0;
3
4 function pickCatalog() {
5     const url = CATALOG_REPLICAS[catIdx % CATALOG_REPLICAS.length];
6     catIdx = (catIdx + 1) % CATALOG_REPLICAS.length;
7     return url;
8 }
9
10 function pickOrder() {
11     const url = ORDER_REPLICAS[ordIdx % ORDER_REPLICAS.length];
12     ordIdx = (ordIdx + 1) % ORDER_REPLICAS.length;
13     return url;
14 }
```

**Note:** To maintain consistency between the primary replica and its backup, all write operations are processed only by the primary replica, which then propagates the updates to the backup replica so that both remain in the same state.

```
1 async function replicateToPeer(op) {
2     if (!IS_PRIMARY) return;
3     try {
4         await fetch(`${PEER_URL}/internal/replicate`, {
5             method: "POST",
6             headers: { "Content-Type": "application/json" },
7             body: JSON.stringify(op),
8         });
9     } catch (e) {
10         console.warn(`[Catalog:${PORT}] replicateToPeer failed:`, e?.
11             message || e);
12     }
13 }
14
15 app.put("/update/price", async (req, res) => {
16     const { id, price } = req.body || {};
17     if (!Number.isInteger(id) || !Number.isFinite(price) || price < 0)
18         return res.status(400).json({ error: "Invalid id or price" });
19
20     await invalidateCache(id);
21     const local = applyWrite({ type: "update_price", id, price });
22     if (!local.ok) return res.status(local.code).json({ error: local.
23         error });
24
25     await replicateToPeer({ type: "update_price", id, price });
26     res.json({ ok: true, primary: IS_PRIMARY, port: PORT });
27 });
```

## 5. Cache Invalidation

To avoid serving outdated data from the frontend cache, a server-push cache invalidation mechanism was implemented. Before carrying out any write operation, the backend services notify the frontend to invalidate the relevant cache entries.

⇒ Before any write, the cache at the frontend service is invalidated.

```
1 // Invalidate endpoint => server-push implementation
```

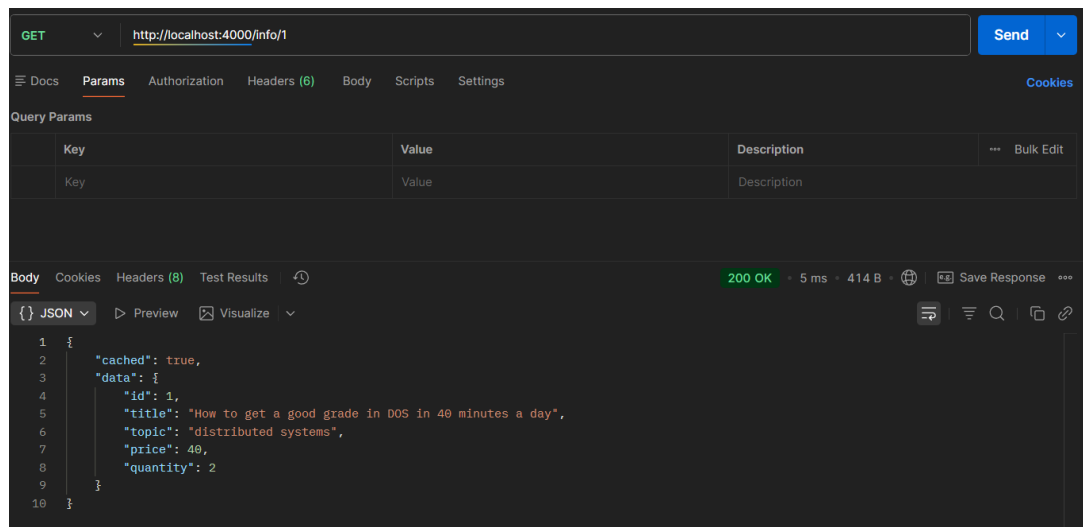
```
2 app.post("/cache/invalidate/:id", (req, res) => {  
3   const id = Number(req.params.id);  
4   if (!Number.isInteger(id) || id <= 0) return res.status(400).json({  
     error: "Invalid id" });  
5  
6   invalidateBook(id);  
7   res.json({ ok: true, invalidated: id });  
8 });
```

**Note:** The Order service communicates with the primary Catalog service, and any update performed by the primary Catalog is propagated to its backup.

## 6. Experiments

Not Cached (miss)  $\Rightarrow$  93ms

Cached (hit)  $\Rightarrow$  5ms



GET <http://localhost:4000/info/1> Send

Docs Params Authorization Headers (6) Body Scripts Settings Cookies

Query Params

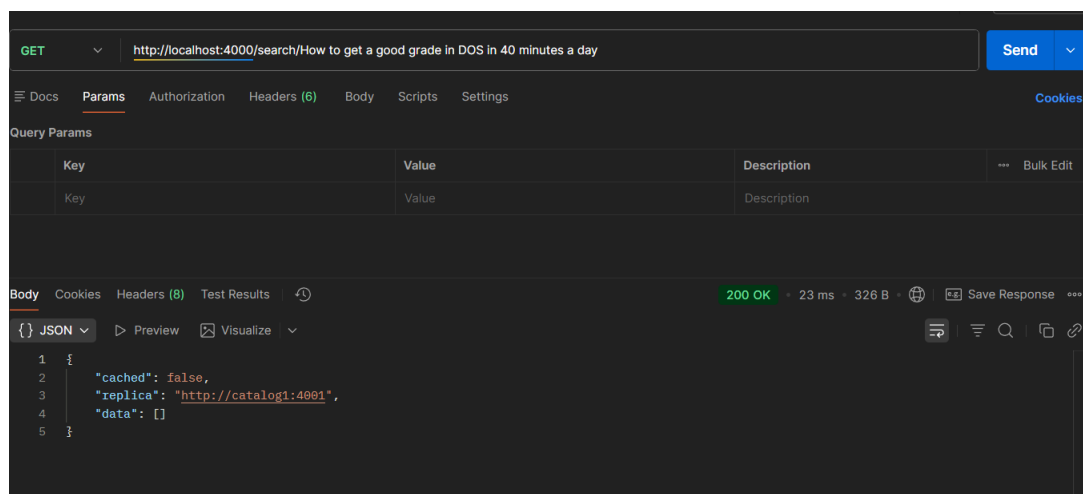
Key	Value	Description	Bulk Edit
Key	Value	Description	

Body Cookies Headers (8) Test Results 200 OK • 5 ms • 414 B Save Response

{ } JSON Preview Visualize

```
1 {
2   "cached": true,
3   "data": {
4     "id": 1,
5     "title": "How to get a good grade in DOS in 40 minutes a day",
6     "topic": "distributed systems",
7     "price": 40,
8     "quantity": 2
9   }
10 }
```

# Not cached (miss)  $\Rightarrow$  23ms



GET <http://localhost:4000/search/How to get a good grade in DOS in 40 minutes a day> Send

Docs Params Authorization Headers (6) Body Scripts Settings Cookies

Query Params

Key	Value	Description	Bulk Edit
Key	Value	Description	

Body Cookies Headers (8) Test Results 200 OK • 23 ms • 326 B Save Response

{ } JSON Preview Visualize

```
1 {
2   "cached": false,
3   "replica": "http://catalog1:4001",
4   "data": []
5 }
```

# Cached (hit)  $\Rightarrow$  4ms

GET

http://localhost:4000/search/How to get a good grade in DOS in 40 minutes a day

Send

Docs

Params

Authorization

Headers (6)

Body

Scripts

Settings

Cookies

Query Params

	Key	Value	Description	...	Bulk Edit
	Key	Value	Description		

Body

Cookies

Headers (8)

Test Results

🕒

200 OK · 4 ms · 292 B · 🌐 | 📄 Save Response · ⋮

{ } JSON

▶ Preview

🖼 Visualize

⌵

1 {

2 |

3 | "cached": true,

4 | "data": []

5 }

## 7. Conclusion

Implementation: introducing caching has greatly decreased the average response time.

Note: The first time a topic is searched, the data is obtained from the catalog service (cache miss). When the same topic is searched again, the data is retrieved directly from the cache, showing a clear improvement in efficiency.

To assess the performance of the in-memory cache, response times were recorded for repeated read requests. The initial request produced a cache miss and needed to communicate with the catalog service, resulting in higher latency. Following requests were served straight from the cache, which markedly reduced the response time.

### **Ports and Services:**

Front: 4000

Catalog 1: 4001

Catalog 2: 4003

Order 1: 4002

Order 2: 4004