# CENG 495

## Cloud Computing

Spring 2024-2025

## Second Assignment

Name Surname: Alkım Doğan
Student ID: 2521482

# Documentation of Each Step

I have tried each step in MacOs ARM architecture M1 chip.

- **Click here** to see the github repo that I have chosen among the projects in the list given in the homework pdf.

- I actually tried a lot of repos and finally was able to find that one which worked fine.

- Then, I just cloned it to my local by using the command below

  `git clone https://github.com/ewolff/microservice-kubernetes.git`

- Now, I followed the quickstart part from the **skaffold** documentation.

- I generated two **.yaml** files for **Ollama**. One of them is for service where the other is for Ollama itself.

- `mvn clean package` is for building the project I have chosen. It is required for running as far as I see from the Readme file in the repository.

**ollama.yaml** is below. I increased the memory since I faced same issues at the first trials regarding the resource requirements related to ollama.

```yaml
apiVersion: apps/v1
kind: Deployment
metadata:
  name: ollama
spec:
  selector:
    matchLabels:
      name: ollama
  template:
    metadata:
      labels:
        name: ollama
    spec:
      containers:
      - name: ollama
        image: ollama/ollama:latest
        ports:
        - name: http
          containerPort: 11434
          protocol: TCP
        resources:
          requests:
            memory: "4Gi"
            cpu: "2000m"
          limits:
            memory: "6Gi"
            cpu: "4000m"
```

**ollama-service.yaml** is below.

```yaml
apiVersion: v1
kind: Service
metadata:
  name: ollama-service
  namespace: default
spec:
  type: LoadBalancer
  selector:
    name: ollama
  ports:
    - port: 11434
      targetPort: 11434
```

I also changed the content of index.html file so that there is a user friendly chatbot in the front-end. The content of the HTML is below. I added bootstrap libraries and very basic HTML and JavaScript.

```html
<!DOCTYPE html>
<html>
<head>
<title>Order Processing</title>

    <link href="https://cdn.jsdelivr.net/npm/bootstrap@5.3.5/dist/css/bootstrap.
        min.css" rel="stylesheet" integrity="sha384-SgOJa3DmI69IUzQ2PVdRZhwQ+dy64/
        BUtbMJw1MZ8t5HZApcHrRKUc4W0kG879m7" crossorigin="anonymous">

</head>
<body>
<div>
  <h1>Order Processing</h1>
  <div class="container">
    <div class="row">
      <div class="col-md-4">
        <a href="/customer/list.html">Customer</a>
      </div>
      <div class="col-md-4">List / add / remove customers</div>
    </div>
    <div class="row">
      <div class="col-md-4">
        <a href="/catalog/list.html">Catalog</a>
      </div>
      <div class="col-md-4">List / add / remove items</div>
    </div>
    <div class="row">
      <div class="col-md-4">
        <a href="/catalog/searchForm.html">Catalog</a>
      </div>
      <div class="col-md-4">Search Items</div>
    </div>
    <div class="row">
      <div class="col-md-4">
        <a href="/order/">Order</a>
      </div>
      <div class="col-md-4">Create an order</div>
    </div>
    <div class="row">
    </div>
  </div>

</div>
<div class="d-flex justify-content-center" style="padding-top: 3rem; margin-top:3
    rem ;">
  <div id="chatbox">
    <h2>Ollama Chatbot</h2>
    <input type="text" id="userInput" placeholder="Type your message...">
    <button id="sendButton">Send</button>
  </div>
</div>

<script>
  const chatbox = document.getElementById('chatbox');
```

```javascript
    const userInput = document.getElementById('userInput');
    const sendButton = document.getElementById('sendButton');

    const CHATBOT_URL = 'http://ollama-service:32761/api/generate';

    sendButton.addEventListener('click', sendMessage);
    userInput.addEventListener('keypress', function (e) {
      if (e.key === 'Enter') {
        sendMessage();
      }
    });

    function appendMessage(sender, text) {
      const messageDiv = document.createElement('div');
      messageDiv.classList.add('message', sender);
      messageDiv.innerText = `${sender === 'user' ? 'You' : 'Bot'}: ${text}`;
      chatbox.appendChild(messageDiv);
      chatbox.scrollTop = chatbox.scrollHeight;
    }

    async function sendMessage() {
      const message = userInput.value.trim();
      if (message === '') return;

      appendMessage('user', message);
      userInput.value = '';

      try {
        const response = await fetch(CHATBOT_URL, {
          method: 'POST',
          headers: {
            'Content-Type': 'application/json'
          },
          body: JSON.stringify({
            model: 'smollm2',
            prompt: message,
            stream: false
          })
        });

        const data = await response.json();
        if (data.response) {
          appendMessage('bot', data.response);
        } else {
          appendMessage('bot', 'Error: ' + JSON.stringify(data));
        }
      } catch (error) {
        console.error('Error:', error);
        appendMessage('bot', 'Failed to connect to the chatbot.');
      }
    }
  }
</script>

<script src="https://cdn.jsdelivr.net/npm/@popperjs/core@2.11.8/dist/umd/popper.
    min.js" integrity="sha384-I7E8VVD/ismYTF4hNIPjVp/Zjvgyol6VFvRkX/vR+Vc4jQkC+
    hVqc2pM8ODewa9r" crossorigin="anonymous"></script>
<script src="https://cdn.jsdelivr.net/npm/bootstrap@5.3.5/dist/js/bootstrap.min.
```

```
        js"  integrity="sha384−VQqxDN0EQCkWoxt/0↩
        vsQvZswzTHUVOImccYmSyhJTp7kGtPed0Qcx8rK9h9YEgx+"  crossorigin="anonymous"≻</↩
        script>

</body>
</html>
```

One thing to mention would be the URL of the ollama in the index.html.   `const CHATBOT_URL = 'http://ollama-service:32761/api/generate';`  is used in the JavaScript code. This URL is used for the Ollama service in the kubernetes. (**NOTE**: `const CHATBOT_URL = 'http://localhost:11434/api/generate'` can also be used for some cases.)

- Now, I run the following commands.

  1. `skaffold init`. (For initializing the skaffold)
  2. `minikube start -memory=7g -cpus=4` (For starting the minikube with memory constrains to allow ollama)
  3. `skaffold config set -global local-cluster true`
  4. `eval $(minikube -p custom docker-env)` (These are directly from documentation of skaffold)
  5. `skaffold dev`. (This command is used for running to begin using Skaffold for continuous development)

Now, we receive the **skaffold.yaml** file automatically generated with the following content.

```yaml
apiVersion: skaffold/v4beta13
kind: Config
metadata:
  name: microservice-kubernetes-demo
build:
  artifacts:
    - image: docker.io/ewolff/microservice-kubernetes-demo-apache
      context: apache
      docker:
        dockerfile: Dockerfile
    - image: docker.io/ewolff/microservice-kubernetes-demo-catalog
      context: microservice-kubernetes-demo-catalog
      docker:
        dockerfile: Dockerfile
    - image: docker.io/ewolff/microservice-kubernetes-demo-customer
      context: microservice-kubernetes-demo-customer
      docker:
        dockerfile: Dockerfile
    - image: docker.io/ewolff/microservice-kubernetes-demo-order
      context: microservice-kubernetes-demo-order
      docker:
        dockerfile: Dockerfile
manifests:
  rawYaml:
    - microservices.yaml
    - ollama-service.yaml
    - ollama.yaml
```

Figure 1: File Paths of The Project. We can see ollama yaml's as well as skaffold.yaml. Index.html can also be seen inside apache server.

After this step, we can also see the some details about the deployment. They actually show that we are on the right path. Below are the screenshots.



Figure 2: Output of kubectl get deployments

Figure 3: Output of kubectl get pods



Figure 4: Output of kubectl get services

Then, It is time to pull ollama to the pod.

- `kubectl -n default exec -ti ollama-6f449bb966-jx6pl - /bin/bash` (For SSH to Ollama service pod).

- `ollama pull smollm2` ( For pulling the lightweight LLM we chose)

The command above is for ssh into pod. The importent issue is that I have made use of Pod Name to ssh into it.



Figure 5: SSH into Pod to download (pull) ollama **smollm2**

Now, I type the command `minikube service apache -url` to see the url. I also visit the url from my browser. The output of the command is **http://127.0.0.1:54077**
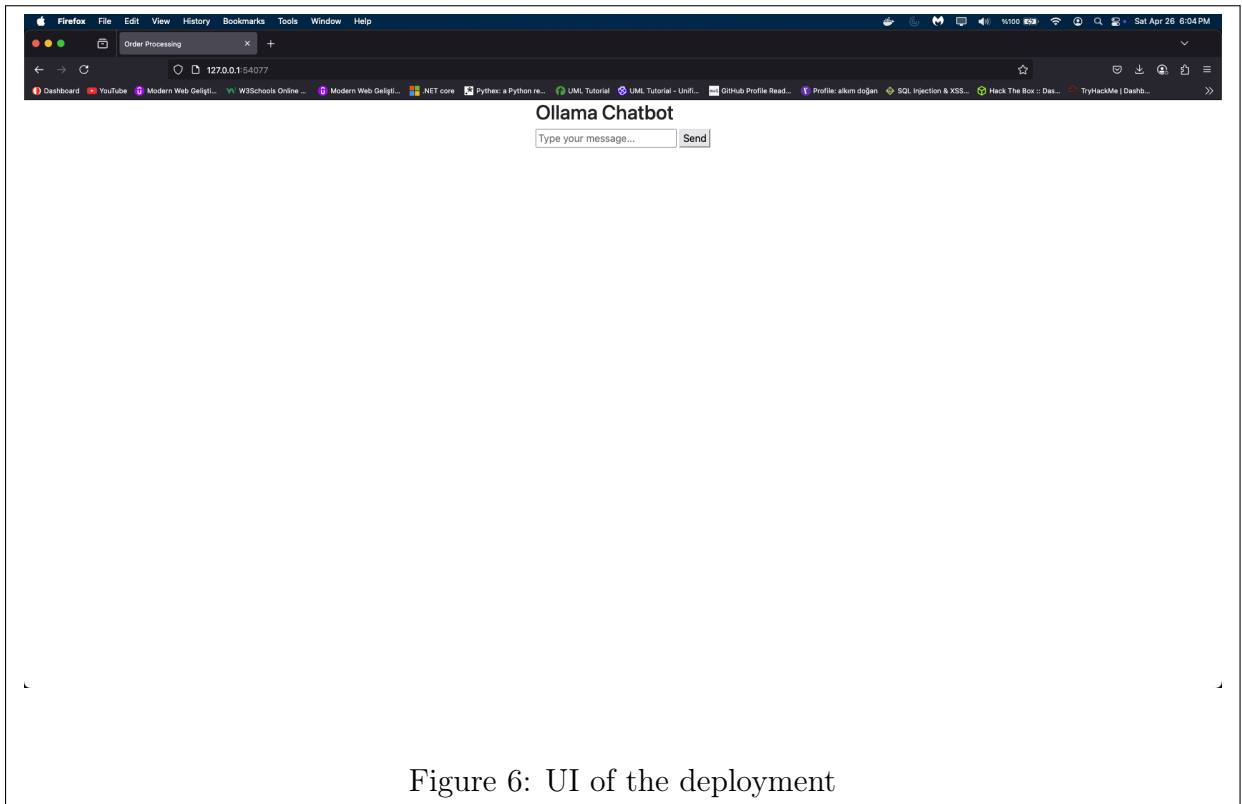
8

Figure 6: UI of the deployment

After the completion of the downloading (pulling) process of ollama, we can directly talk to chatbot as follows.
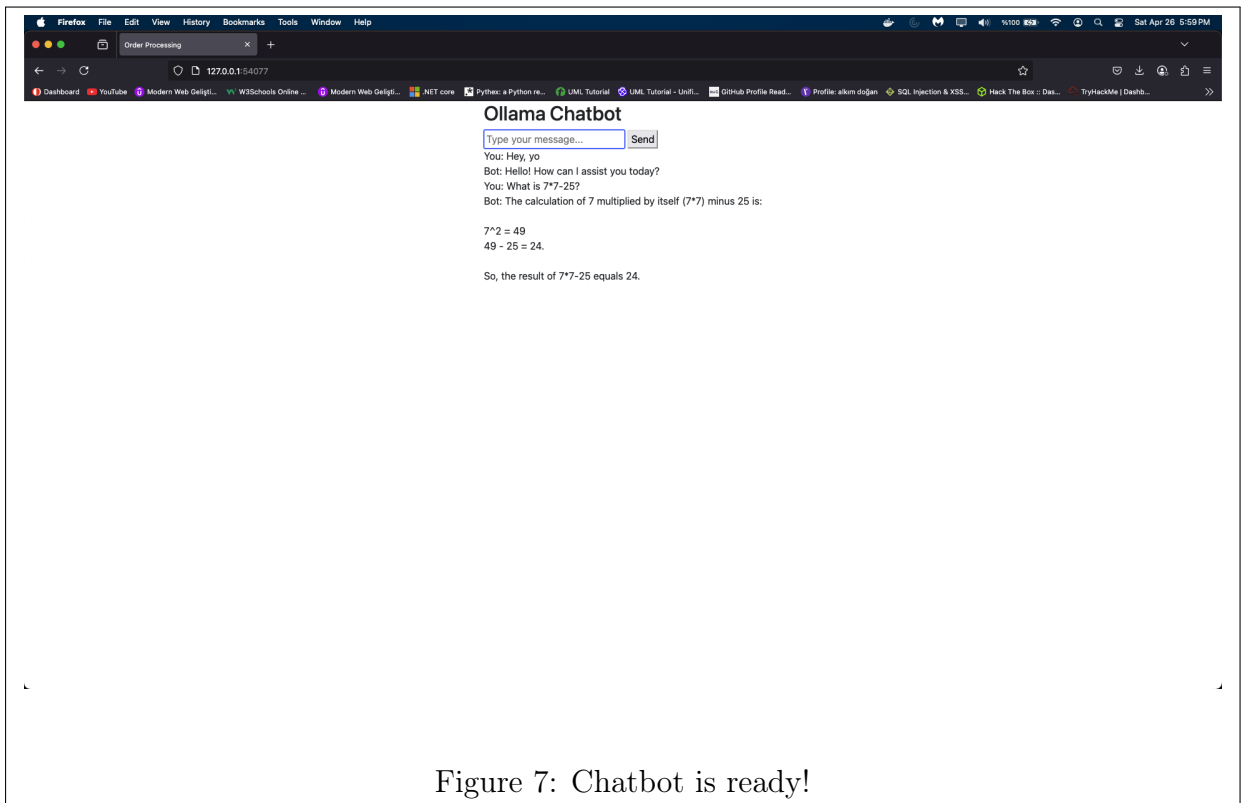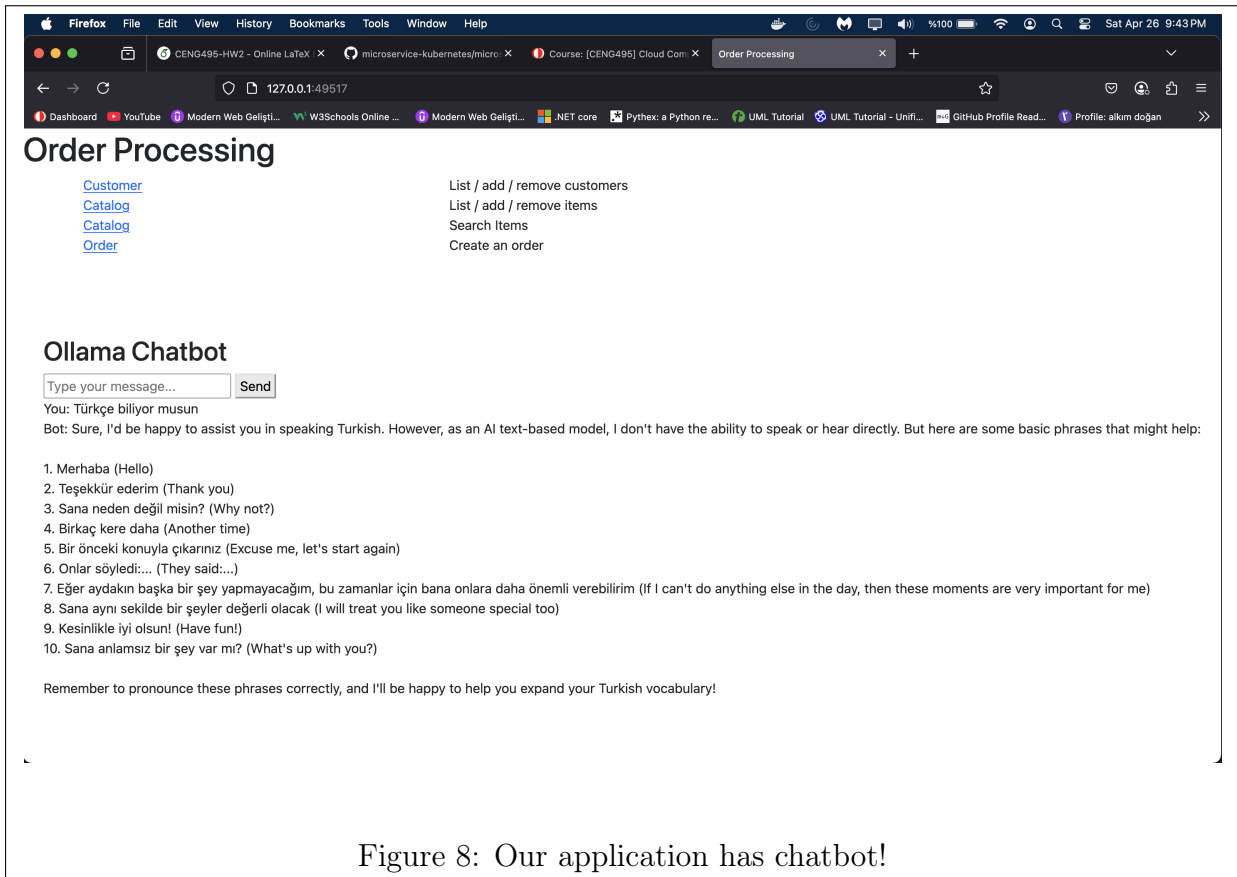


Figure 7: Chatbot is ready!

Figure 8: Our application has chatbot!

In the previous image, I had changed the index.html so we can only see chatbot. I re-run the project and changed index.html so that both the application I cloned and ollama chat are available in the front end html file.

**Additional Commands**

- `kubectl port-forward pod/<pod_name> 11434:11434` (This is used for port forwarding of given pod)

- `kubectl apply -f <file_name>.yaml` (For creating the service given the file name of the yaml )