

DevOps Implementation of a Containerized Web Scrapping Application

Overview:

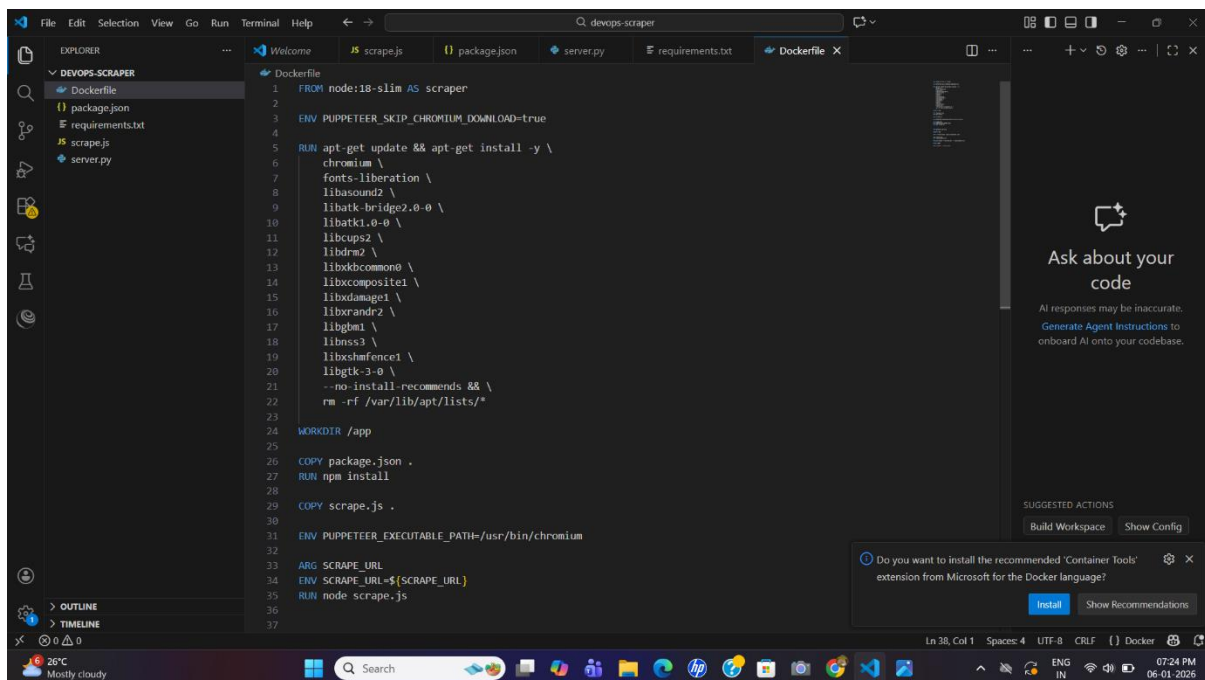
The project creates and executes a containerized web scrapping application using Docker. Basic website data is scraped using a Node.js service with Puppeteer, and the extracted output is served through a Python flask. Docker multi-stage builds are used to separate the scraping and server layers, resulting in a lightweight and efficient final image. The application demonstrates the practical use of Docker for application packaging and deployment, as it runs successfully inside a container and is accessible using a local web interface.

Tools Used:

- **Docker & Docker Desktop** – For containerization and multi-stage builds
- **Node.js** – To run the web scraping service
- **Puppeteer** – For automated web scraping using a headless browser
- **Python** – To develop the backend service
- **Flask** – To serve scraped data through a REST API
- **Chromium** – Headless browser used by Puppeteer
- **Visual Studio Code** – Development environment
- **Web Browser (Chrome)** – To access the API output

Assignment Explanation:

Dockerfile – Stage 1: Web Scrapping (Node.js & Puppeteer)



```
1 FROM node:18-slim AS scraper
2
3 ENV PUPPETEER_SKIP_CHROMIUM_DOWNLOAD=true
4
5 RUN apt-get update && apt-get install -y \
6     chromium \
7     fonts-liberation \
8     libasound2 \
9     libatk-bridge2.0-0 \
10    libatk1.0-0 \
11    libcupst2 \
12    libdrm2 \
13    libfontconfig \
14    libfreetype6 \
15    libgdk-pixbuf2.0-0 \
16    libglib2.0-0 \
17    libgtk-3-0 \
18    libnss3 \
19    libpango1.0-0 \
20    libx11-xcb1 \
21    libxcomposite1 \
22    libxrandr2 \
23    --no-install-recommends && \
24    rm -rf /var/lib/apt/lists/*
25
26 WORKDIR /app
27
28 COPY package.json .
29 RUN npm install
30
31 COPY scraper.js .
32
33 ENV PUPPETEER_EXECUTABLE_PATH=/usr/bin/chromium
34
35 ARG SCRAPE_URL
36 ENV SCRAPE_URL=${SCRAPE_URL}
37 RUN node scraper.js
```

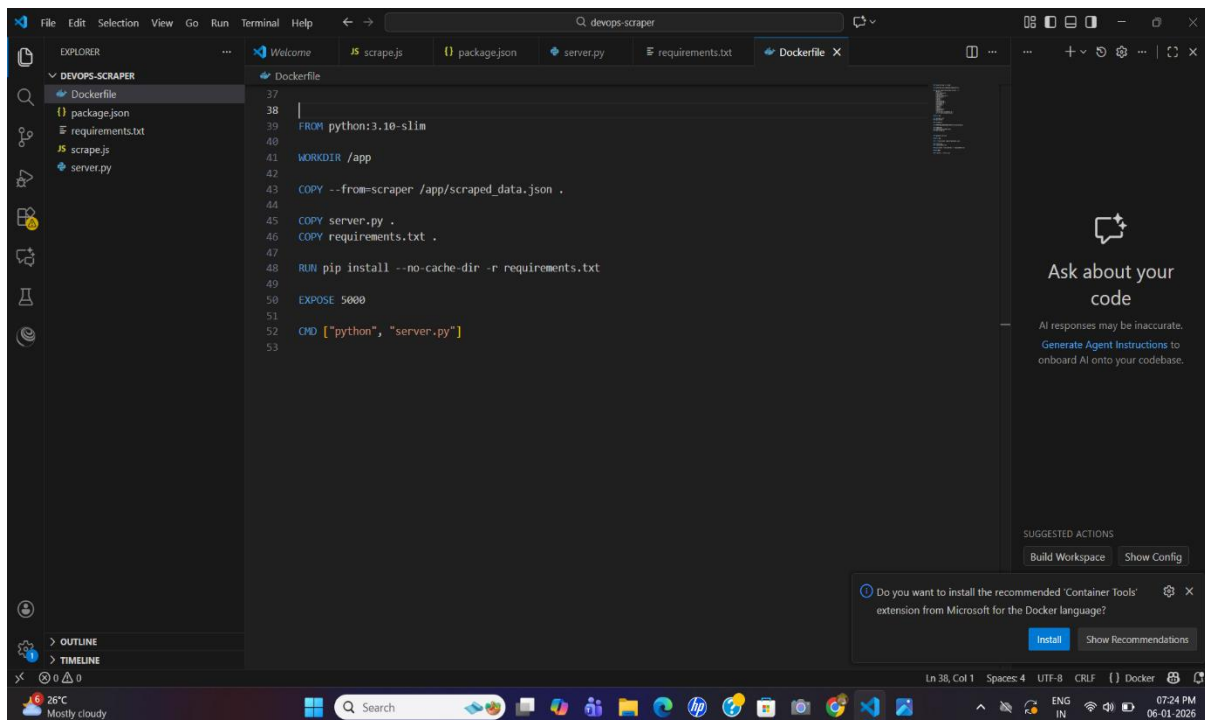
This stage of the Dockerfile is responsible for performing the web scraping operation using Node.js and Puppeteer. The base image node:18-slim is used to keep the image lightweight while supporting Node.js execution.

Puppeteer requires a browser environment, so Chromium and its required system libraries are installed using apt-get. The environment variable is set to skip the default Chromium download and instead use the system-installed browser.

The working directory is set to /app, and Node.js dependencies are installed using npm install. The scraping script (scrape.js) is then copied into the container. An environment variable is used to pass the target website URL dynamically.

Finally, the Node.js script is executed, which launches a headless browser, scrapes the webpage title and first heading, and saves the output as a JSON file. This output is later used by the second stage of the Docker build.

Dockerfile – Stage 2: Flask Server (Python)



```
37 |
38 | FROM python:3.10-slim
39 |
40 |
41 | WORKDIR /app
42 |
43 | COPY --from=scrapper /app/scraped_data.json .
44 |
45 | COPY server.py .
46 | COPY requirements.txt .
47 |
48 | RUN pip install --no-cache-dir -r requirements.txt
49 |
50 | EXPOSE 5000
51 |
52 | CMD ["python", "server.py"]
53 |
```

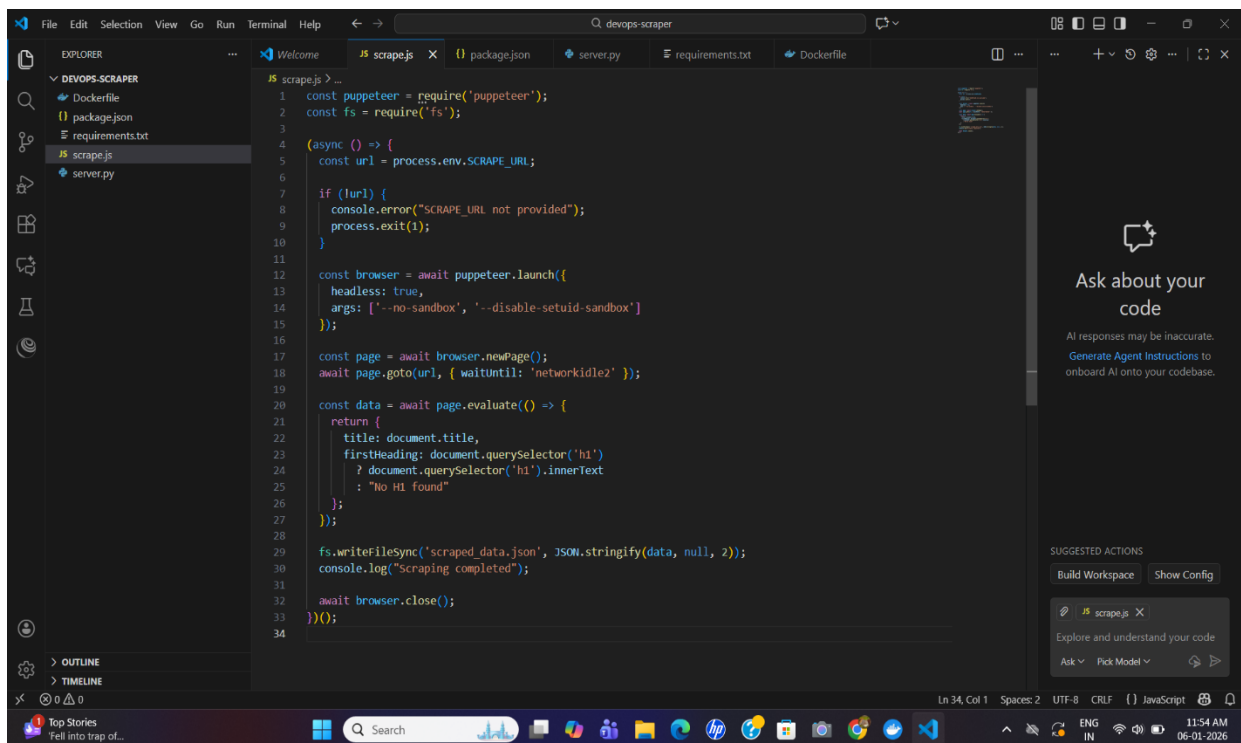
This stage of the Dockerfile is responsible for serving the scraped data using a Python Flask application. The base image python:3.10-slim is used to keep the final container lightweight while supporting Python execution.

The working directory is set to /app. The scraped JSON file generated in the first (Node.js) stage is copied into this stage using a multi-stage build approach. This avoids carrying unnecessary Node.js and browser dependencies into the final image.

The Flask server file (server.py) and the dependency file (requirements.txt) are then copied into the container. Required Python libraries are installed using pip, ensuring only essential packages are included.

Port 5000 is exposed to allow external access to the Flask application. Finally, the container starts the Flask server using the Python command, making the scraped website data accessible through a local web interface.

Web Scraping Script (scrape.js)



```
1 const puppeteer = require('puppeteer');
2 const fs = require('fs');
3
4 (async () => {
5   const url = process.env.SCRAPE_URL;
6
7   if (!url) {
8     console.error('SCRAPE_URL not provided');
9     process.exit(1);
10  }
11
12  const browser = await puppeteer.launch({
13    headless: true,
14    args: ['--no-sandbox', '--disable-setuid-sandbox']
15  });
16
17  const page = await browser.newPage();
18  await page.goto(url, { waitUntil: 'networkidle2' });
19
20  const data = await page.evaluate(() => {
21    return {
22      title: document.title,
23      firstHeading: document.querySelector('h1')
24        ? document.querySelector('h1').innerText
25        : 'No H1 found'
26    };
27  });
28
29  fs.writeFileSync('scraped_data.json', JSON.stringify(data, null, 2));
30  console.log('Scraping completed');
31
32  await browser.close();
33 })();
34
```

This JavaScript file is responsible for scraping data from a given website using the Puppeteer library. Puppeteer is a Node.js tool that controls a headless Chromium browser to interact with web pages programmatically.

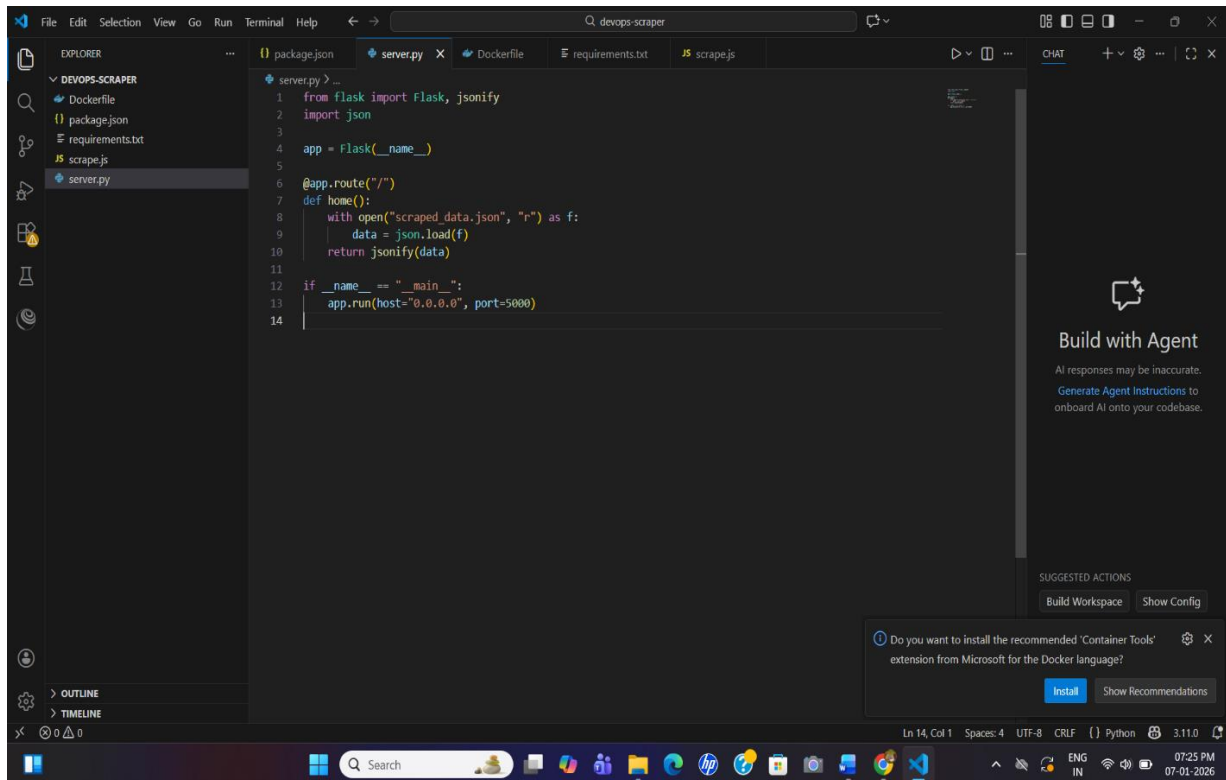
The script first reads the target website URL from an environment variable called `SCRAPE_URL`. This makes the script flexible, as different websites can be scraped without changing the code. If the URL is not provided, the program stops execution and displays an error message.

A headless browser is then launched with sandbox restrictions disabled to ensure smooth execution inside a Docker container. A new page is opened and the browser navigates to the given URL, waiting until all network requests are completed.

Once the page is fully loaded, the script extracts two pieces of information: the page title and the first heading (`<h1>`). If no heading is found, a default message is returned. The scraped data is saved into a JSON file named `scraped_data.json`, which is later used by the Python Flask server.

Finally, the browser is closed after successful scraping to free system resources.

Flask Server Script (server.py)



```
server.py > ...
1 from flask import Flask, jsonify
2 import json
3
4 app = Flask(__name__)
5
6 @app.route("/")
7 def home():
8     with open("scraped_data.json", "r") as f:
9         data = json.load(f)
10        return jsonify(data)
11
12 if __name__ == "__main__":
13     app.run(host="0.0.0.0", port=5000)
14
```

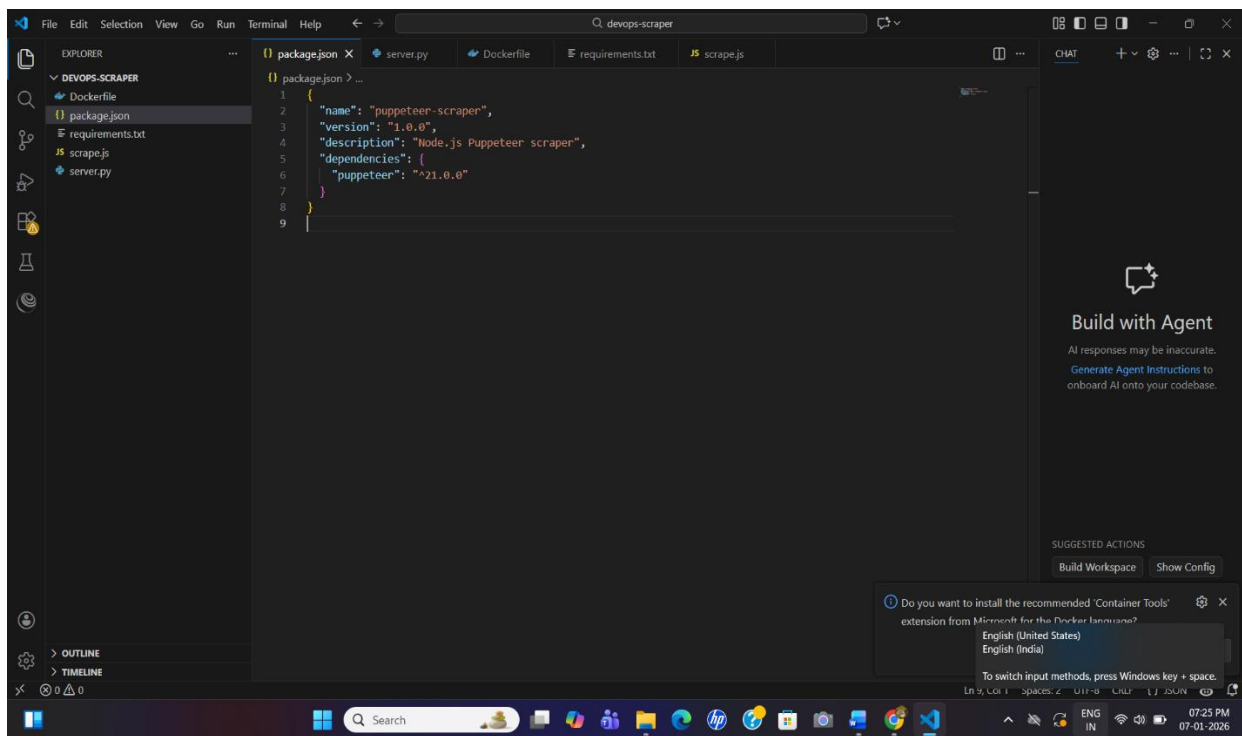
This Python file implements a simple web server using the Flask framework. Its main purpose is to serve the data that was scraped by the Node.js Puppeteer script and stored in a JSON file.

The Flask application starts by importing the required modules, including Flask for creating the web server and json for reading the scraped data file. When the server starts, it looks for the scraped_data.json file in the application directory.

A single route (/) is defined in the application. When a user accesses this route through a web browser, the server reads the contents of the JSON file and returns the data as a JSON response. This allows users to easily view the scraped website title and heading through a web interface.

The server runs on port 5000 and is configured to listen on all network interfaces, making it accessible from outside the Docker container. This setup demonstrates how scraped data can be exposed using a lightweight Python API.

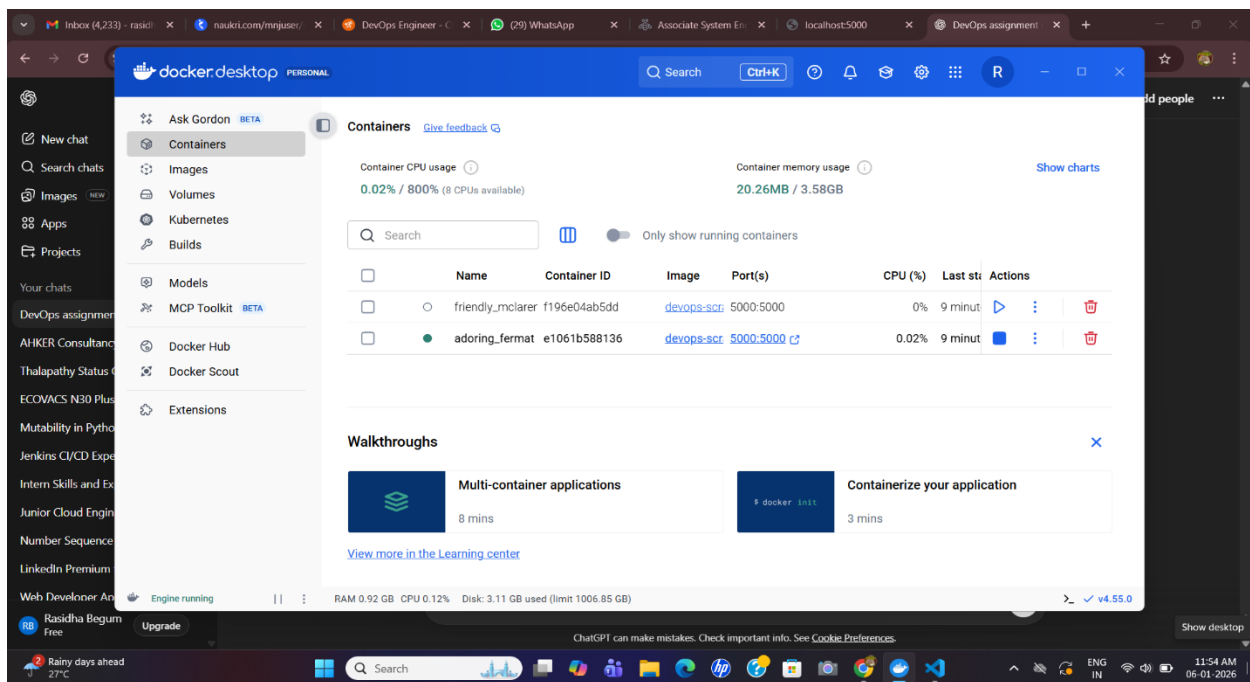
Node.js Configuration File (package.json)



This package.json file defines the configuration and dependencies required for the Node.js scraping component of the project. It specifies basic project details such as the project name, version, and a short description indicating that the application is a Puppeteer-based web scraper.

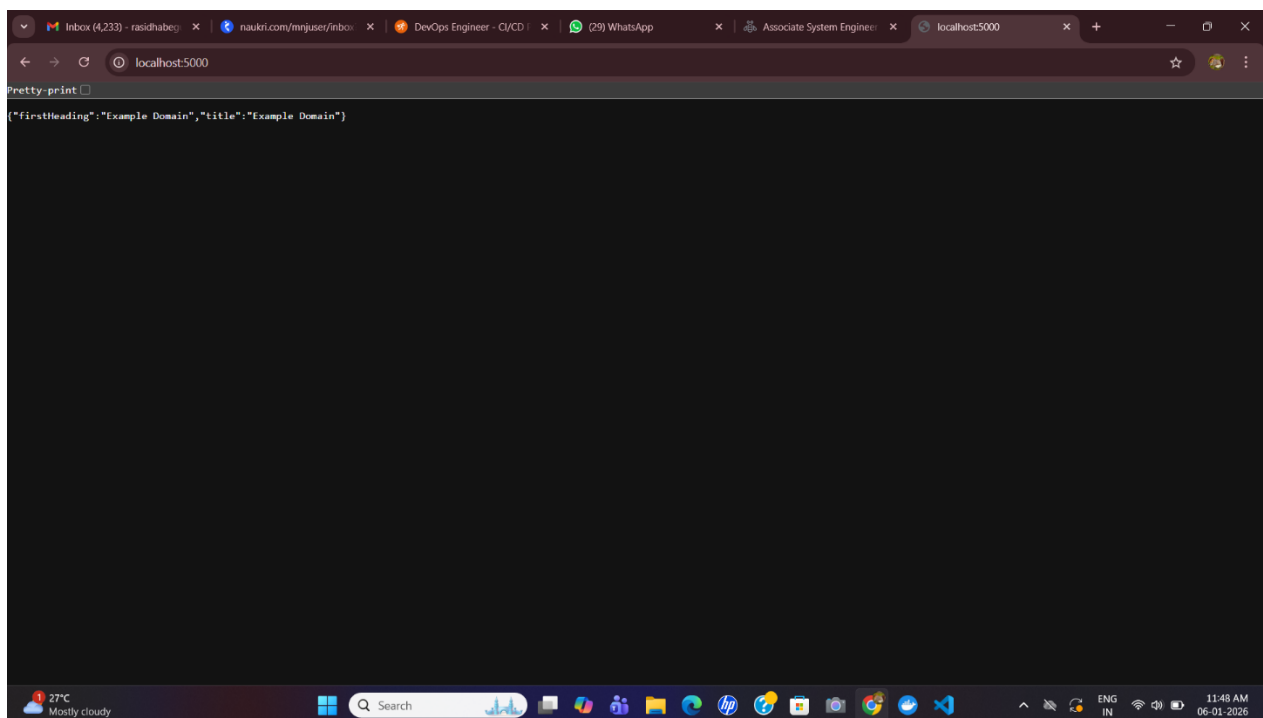
The most important part of this file is the dependencies section, where Puppeteer is listed. Puppeteer is a Node.js library used to control a headless Chrome browser, which enables automated web page loading and data extraction. During the Docker build process, this file is used to install all required Node.js dependencies inside the container.

Step 2: Docker Container Execution



This screen shows the running Docker containers for the project in Docker Desktop. It confirms that the application image is successfully deployed as a container and is actively running, with port **5000** mapped to the host system for accessing the Flask web service through a browser.

Step 3: Final Output Verification



This screenshot shows the final output of the application accessed through **http://localhost:5000**. The Flask server successfully displays the scraped website data in JSON format, including the page title and first heading, confirming that the containerized scraping and API services are working correctly.

Challenges Faced:

- Docker image build failed due to network DNS issues.
- Puppeteer failed initially because Chromium was not installed.
- Port 5000 was not accessible at first.

These issues were resolved by reinstalling Docker, configuring DNS, and rebuilding the image.