# usrbinkat / README.md

Last active 3 days ago • Report abuse

---

<> **Code**      -O- **Revisions**   3      ☆ **Stars**   48      ⑂ **Forks**   12

---

Ollama + Open-Webui + Nvidia/CUDA + Docker + docker-compose

<> **README.md**



UPDATE: This is tested and working on both Linux and Windows 11 used for LlaMa & DeepSeek

Here's a sample `README.md` file written by Llama3.2 using this docker-compose.yaml file that explains the purpose and usage of the Docker Compose configuration:

**ollama-portal**

A multi-container Docker application for serving OLLAMA API.

**Overview**

This repository provides a Docker Compose configuration for running two containers: `open-webui` and `ollama`. The `open-webui` container serves a web interface that interacts with the `ollama` container, which provides an API or service. This setup is designed to work together seamlessly, allowing users to access OLLAMA's functionality through a user-friendly web interface.

### Architecture

The application consists of two main components:

- **OLLAMA**: A GPU-accelerated neural network inference service that provides a RESTful API for serving models.
- **Open-WebUI**: A web-based interface for interacting with the OLLAMA API, providing a simple and intuitive way to deploy and manage models.

### Docker Compose Configuration

The Docker Compose configuration file (`docker-compose.yaml`) defines several key settings:

- **Services**: The application consists of two services: `open-webui` and `ollama`. Each service is defined with its own set of environment variables, volumes, and ports.
- **Environment Variables**:

- `MODEL_DOWNLOAD_DIR` : Specifies the directory for storing downloaded models.
- `OLLAMA_API_BASE_URL` : Sets the base URL for the OLLAMA API.
- `LOG_LEVEL` : Configures the log level for both containers.

- **Volumes**: The application mounts several volumes to share data between containers. These include:

- `data` : For storing user input and model artifacts.
- `models` : For accessing pre-trained models.
- `ollama` : For storing application-specific data.

### Container Configuration

The Docker Compose configuration defines the following container configurations:

- **OLLAMA Container**:

- Uses the official OLLAMA image (`ollama/ollama:latest`).
- Specifies NVIDIA GPU acceleration using the `runtime: nvidia` option.
- Configures the container to use all available GPUs in the cluster.

- **Open-WebUI Container**:

- Uses the official Open-WebUI image ( `ghcr.io/open-webui/open-webui:main` ).
- Specifies environment variables for model download directories and OLLAMA API URLs.

### Networking

The application uses a single network ( `ollama-net` ) that connects both containers. This allows them to communicate with each other seamlessly.

### Running in Production

To run this application in production, you'll need to:

- Set up your OLLAMA API on the `ollama` container.
- Configure the `open-webui` container to connect to your OLLAMA API.
- Mount necessary volumes and adjust configuration variables as needed.

### Troubleshooting

If you encounter issues while running this application, please refer to the [Docker Compose troubleshooting guide](https://...) for assistance.

### Security Considerations

This application uses the following security measures:

- **Model signing**: The OLLAMA API verifies model signatures using a digital certificate.
- **Input validation**: The Open-WebUI container validates user input to prevent injection attacks.
- **Encryption**: Data exchanged between containers is encrypted using SSL/TLS.
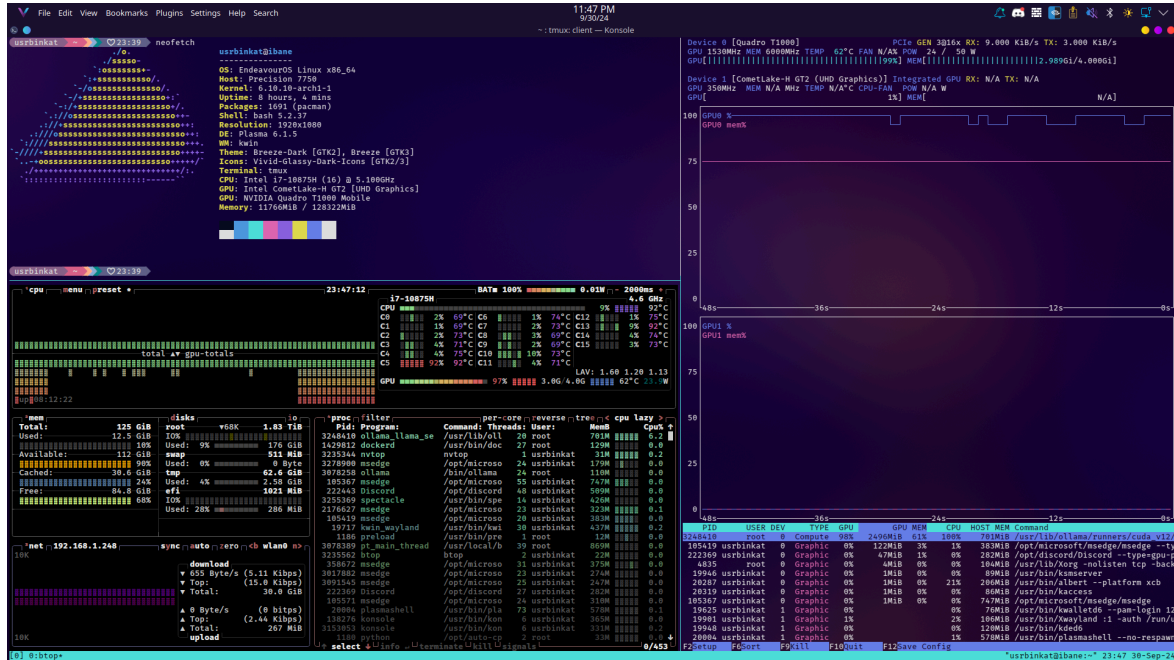
### Performance Optimization

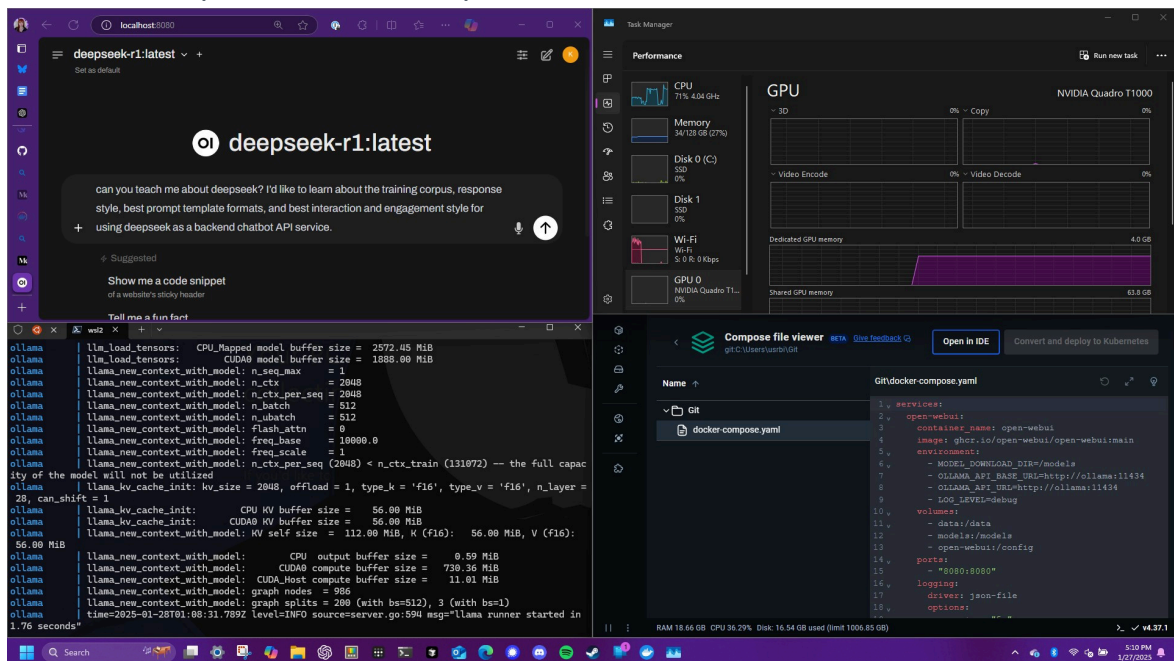To optimize performance, consider the following:

- **Model caching**: Use a caching layer (e.g., Redis) to store frequently accessed models.
- **Container orchestration**: Use a container orchestration tool (e.g., Kubernetes) to manage and scale your containers.
- **GPU acceleration**: Configure multiple GPUs on your system for optimal performance.

This enhanced README file provides more in-depth technical explanations, including architecture, Docker Compose configuration, container configurations, networking, security considerations, and performance optimization. If you have any further questions or concerns, feel free to open a discussion on our GitHub page!

### Arch Linux (Elementary OS)



### Windows 11 (circa: 2025/01/27)



<> **docker-compose.yaml**

```
1   # https://gist.githubusercontent.com/usrbinkat/de44facc683f954bf0cca6c87e2f9f88/raw/040
2   services:
3     open-webui:
4       container_name: open-webui
```

```yaml
 5        image: ghcr.io/open-webui/open-webui:main
 6        environment:
 7          - MODEL_DOWNLOAD_DIR=/models
 8          - OLLAMA_API_BASE_URL=http://ollama:11434
 9          - OLLAMA_API_URL=http://ollama:11434
10          - LOG_LEVEL=debug
11          - WEBUI_SECRET_KEY=your_secret_key_here  # Add this to prevent logouts after upda
12        volumes:
13          - data:/data
14          - models:/models
15          - open-webui:/app/backend/data  # Corrected path based on documentation
16        ports:
17          - "8080:8080"
18        logging:
19          driver: json-file
20          options:
21            max-size: "5m"
22            max-file: "2"
23        depends_on:
24          - ollama
25        extra_hosts:
26          - "host.docker.internal:host-gateway"
27        networks:
28          - ollama-net
29        restart: unless-stopped
30
31      ollama:
32        container_name: ollama
33        image: ollama/ollama:latest
34        runtime: nvidia
35        environment:
36          - NVIDIA_VISIBLE_DEVICES=all
37          - NVIDIA_DRIVER_CAPABILITIES=compute,utility
38          - CUDA_VISIBLE_DEVICES=0
39          - LOG_LEVEL=debug
40        deploy:
41          resources:
42            reservations:
43              devices:
44                - driver: nvidia
45                  capabilities: [gpu]
46                  count: all
47        volumes:
48          - ollama:/root/.ollama
49          - models:/models
50        ports:
51          - "11434:11434"
52        logging:
53          driver: json-file
54          options:
55            max-size: "5m"
56            max-file: "2"
```

```
57        networks:
58          - ollama-net
59        restart: unless-stopped
60
61      watchtower:
62        image: containrrr/watchtower
63        container_name: watchtower
64        volumes:
65          - /var/run/docker.sock:/var/run/docker.sock
66        command: --interval 300 open-webui  # Check for updates every 5 minutes
67        depends_on:
68          - open-webui
69        networks:
70          - ollama-net
71        restart: unless-stopped
72
73    volumes:
74      data:
75      models:
76      ollama:
77      open-webui:
78    networks:
79      ollama-net:
80        driver: bridge
```

**pdoyle12** commented on Mar 13

This is fantastic; thank you. This really ought to get a reference from the OpenWebUI docs, since it's a lot less hassle than fiddling with the containers manually.

I'm not certain this was the issue, but I wasn't able to get this to work until I renamed OLLAMA_API_BASE_URL to OLLAMA_BASE_URL. According to the OpenWebUI documentation, OLLAMA_API_BASE_URL is deprecated. Before that, although the containers did start correctly, openwebui didn't seem to respect the Ollama URL.

**RockportTigger** commented on Jun 16

> This is fantastic; thank you. This really ought to get a reference from the OpenWebUI docs, since it's a lot less hassle than fiddling with the containers manually.

>> other interesting options an possible addition to this! here -->

Optimize Open WebUI: Three practical extensions for a better user experience

**jonndoe47pp** commented on Jul 1

Right got it working (nearly) with a few caveats. After installing the nvidia controller toolkit had to:
sudo nvidia-ctk runtime configure --runtime=docker
Last issue, cant get ollama to find any llms. So think the bridge is not working. HELP! I am a newbie.
So, read i need to add the
--network=host flag, but thjats in a straight docker statement not docker compose , help total
newbie....

**ghost** commented on Jul 13

> Right got it working (nearly) with a few caveats. After installing the nvidia controller toolkit had
> to: sudo nvidia-ctk runtime configure --runtime=docker Last issue, cant get ollama to find any
> llms. So think the bridge is not working. HELP! I am a newbie. So, read i need to add the --
> network=host flag, but thjats in a straight docker statement not docker compose , help total
> newbie....

this worked for me...had same issue, then after sudo systemctl restart docker I retried deploying
stack and it executed great!

The webUI is not super easy to figure out how to add the models. this link helped me:

https://docs.openwebui.com/getting-started/quick-start/starting-with-ollama

use this link to see whats available. gemma3 qwen3 llama3.1 to name a few. then enter the model
name into download model and it will pull it.