# DeepComputing

# DC-ROMA RISC-V Mainboard II AI Model User Guide (Ubuntu AI Image)

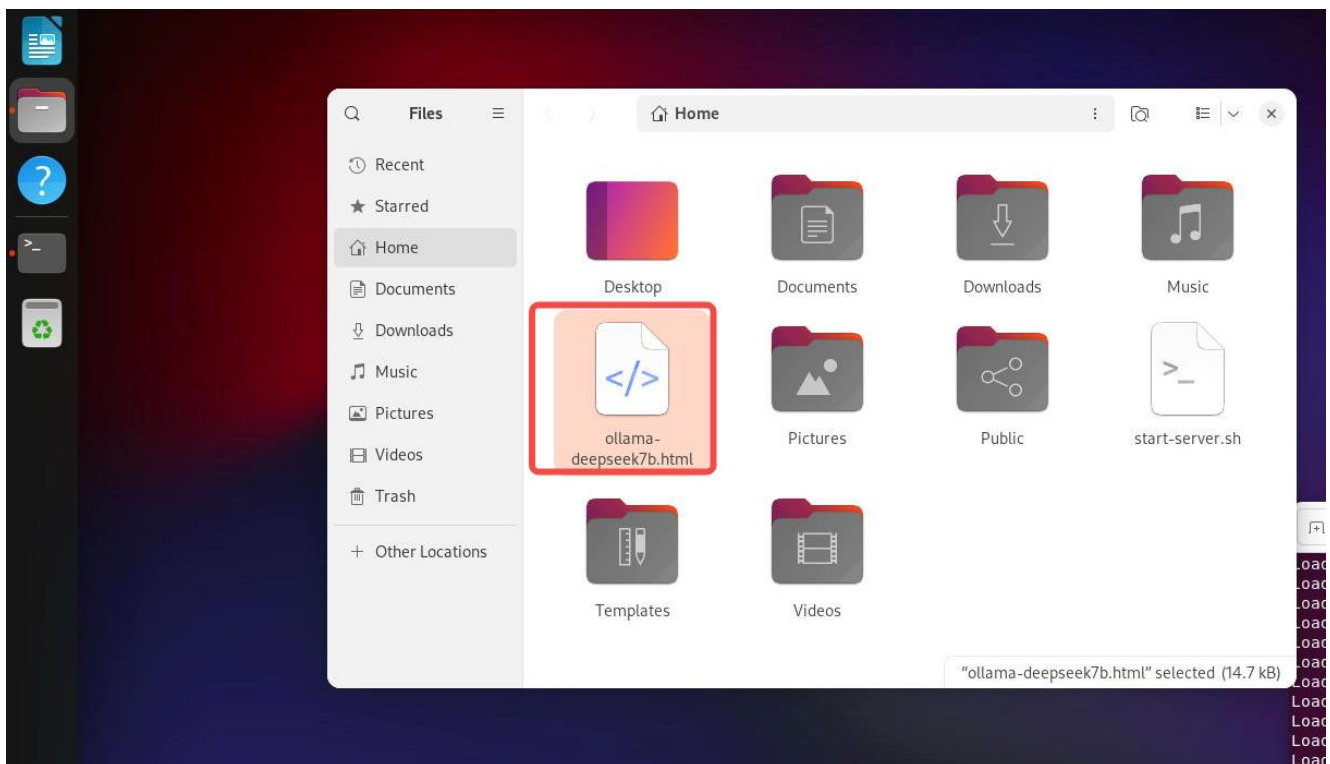| Document Properties | |
|---|---|
| Document version | V1.0 |
| Latest release date | Oct 28, 2025 |
| Applicable product model | DC-ROMA RISC-V Mainboard II for Framework Laptop 13 |
| Operating system version | Ubuntu 24.04 |
| Document target audience | Users holding the DC-ROMA RISC-V Mainboard II series products |
| Document overview | This document applies to users of the factory-installed image system or those who download the image from: http://120.92.155.32:8082/artifactory/virtOS/fml13v03-eswin/15019-ubuntu-24.04-desktop-grub-sdcard-AI.zip |

# Catalog

# Ollama Containerized Deepseek-R1-7B Usage Guide

**The AI PC comes pre-installed with a containerized local DeepSeek 7B model via Ollama and includes a startup script. Simply run the script to use the model. For an enhanced experience, a visual frontend is included.**
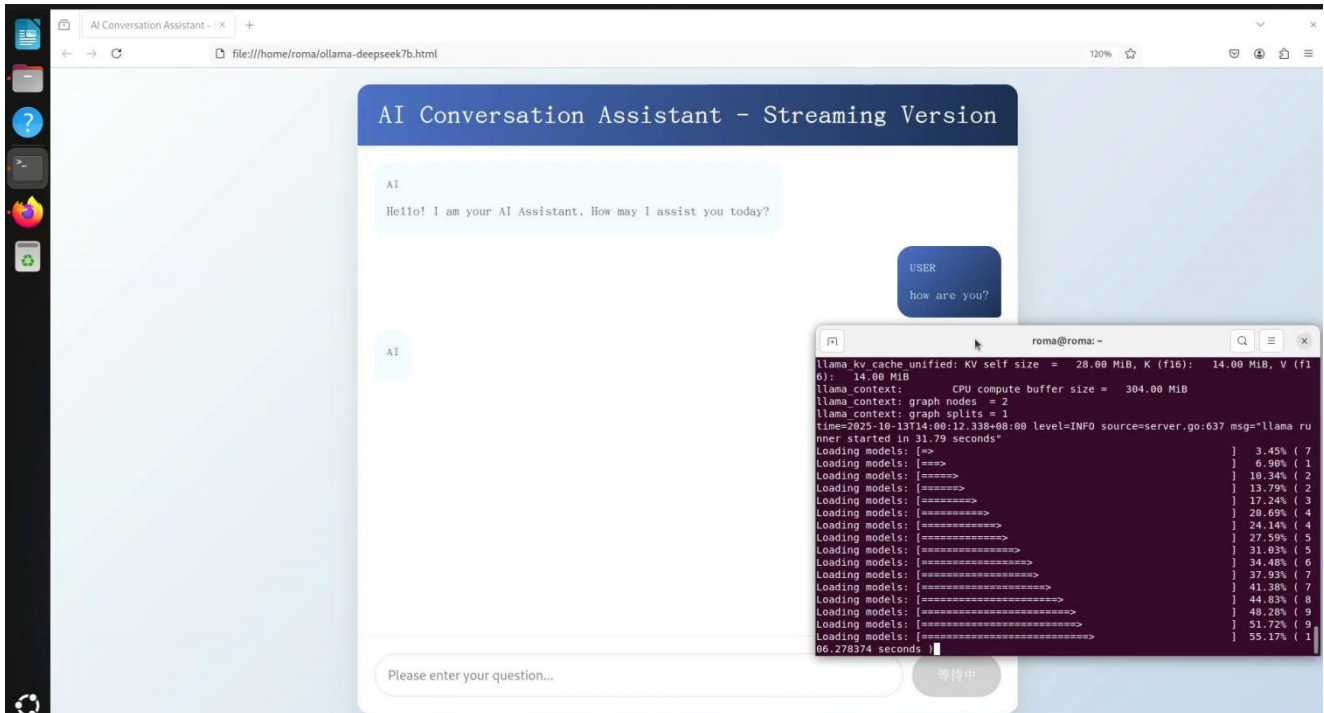
**1、 Open Terminal and run:**
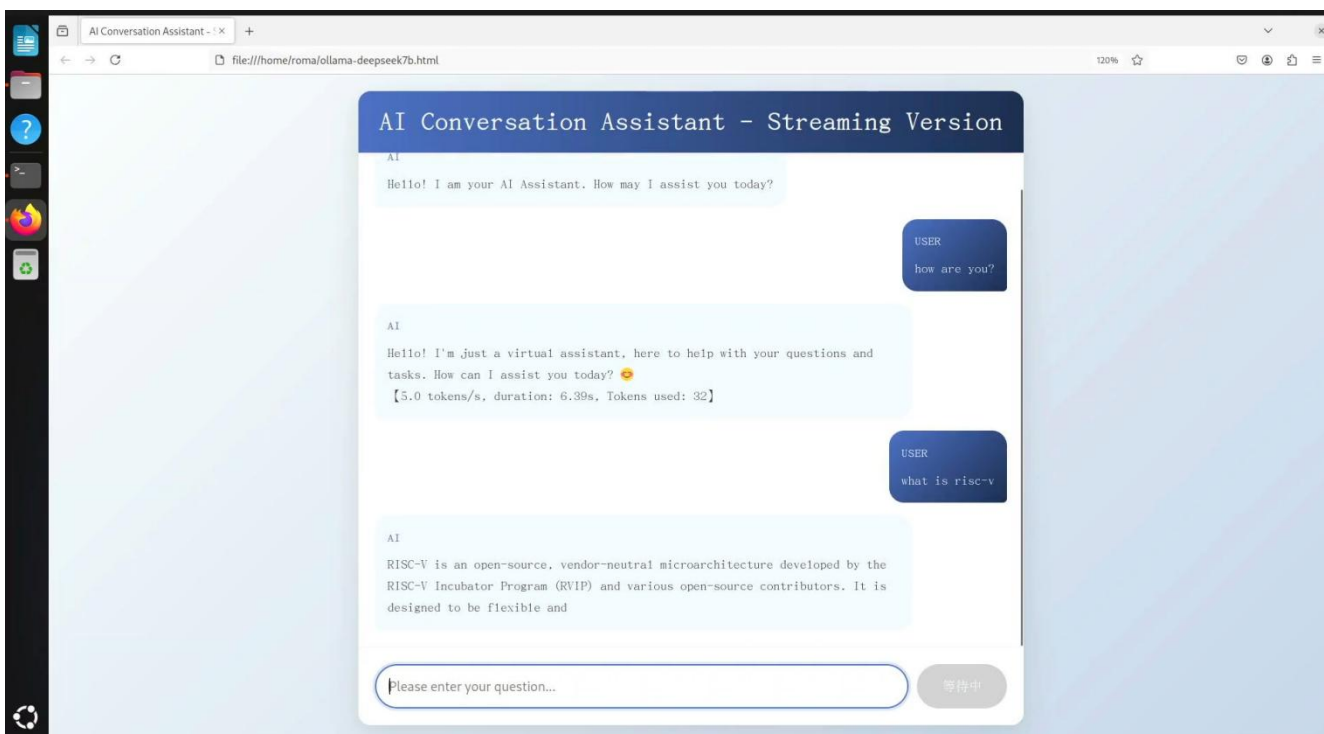
```
sudo  ./start-server.sh
```

**2、 Double-click <u>ollama-deepseek7b.html</u> to open the visual interface.**
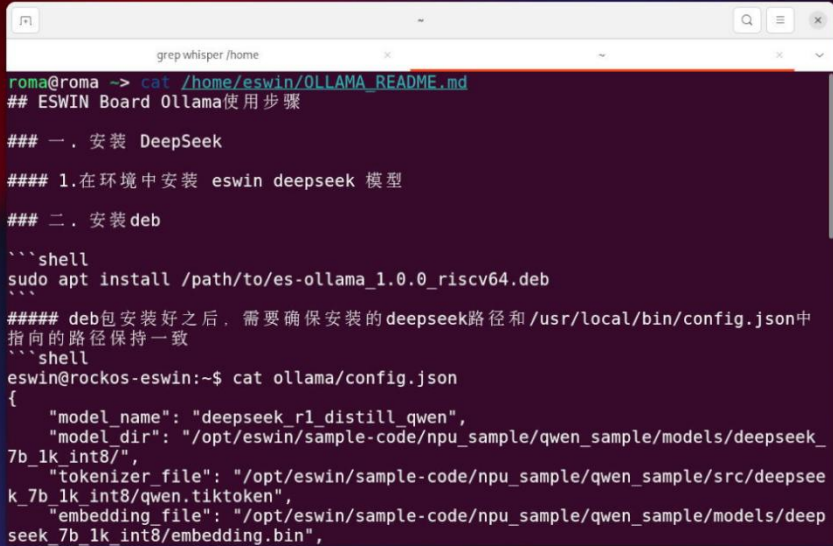


**3、 Ask questions in the dialog box. Note: Initial query loads the model (~3 minutes).**

## 4、 Sample Execution

**5、You can view the usage steps in the /home/eswin/OLLAMA_README.mddirectory.**



# Terminal Execution: Deepseek-7B with Dual-Gen NPU Acceleration

**1、 Launch Model via Terminal:**

```
sudo /opt/eswin/sample-code/npu_sample/qwen_sample/bin/es_qwen2
/opt/eswin/sample-
code/npu_sample/qwen_sample/src/deepseek_7b_1k_int8_peer/config.json
```

**2、 Select the interaction pattern according to the prompt**

```
[root@fedora-riscv roma]# sudo /opt/eswin/sample-code/npu_sample/qwen_sample/bin/es_qwen2 /opt/eswin/sample-code/npu_sample/qwen_sample/src/deepseek_7b_1k_int8_peer/config.json
[E][ES_MEM]  open_malloc_dev:  439 open /dev/malloc_dmabuf failed!
Loading models: [===============================================] 100.00% ( 121.316900 seconds )
------------------------------------------------------------------------------
0: Role setting: 你是一个智能助理.
------------------------------------------------------------------------------
1: 介绍一下大语言模型
2: The quantum computers
3: Humans and robots coexist
4: Customized prompts
------------------------------------------------------------------------------
[YOU]: 4
[YOU]: how are you?
[Qwen2]: <think>

</think>

Hello! I'm just a virtual assistant, here to help you with whatever you need. How can I assist you today?
--------------------------------------------------------
Throughput: 10.3288tokens/s
--------------------------------------------------------
```

# Terminal Execution: Whisper Model

**1、Open Terminal and run the command to view the model usage guide:**

cat  /home/eswin/whisper/WHISPER_README.md


# es_whisper_deb using

\`\`\`shell

# using deb

sudo apt install ./es-whisper_1.0.0_riscv64.deb

# run whisper-cli

/usr/bin/whisper-cli -f /opt/eswin/data/npu/whisper_models/audio/jfk.wav

\`\`\`

**2、 Enter the command to launch the model, which will run a demo converting speech to text.**

```
sudo /usr/bin/whisper-cli -f
/opt/eswin/data/npu/whisper_models/audio/jfk.wav
```