

# How to run Qwen 3.5 locally on a single GPU

Qwen 3.5 is Alibaba's latest Qwen model line, built upon the powerful performance of previous Qwen models in inference, programming, and multimodal tasks.

Qwen 3.5 is Alibaba's latest Qwen model line, built upon the powerful performance of previous Qwen models in inference, programming, and multimodal tasks.

Independent benchmark tests show that the Qwen 3.5-397B-A17B model scores highly in widely used benchmarks such as LiveCodeBench and AIME26, often outperforming leading models like GPT-5.2 and Claude Opus 4.5 in most evaluated categories, and delivering significantly higher throughput compared to previous Qwen generations.

## Hardware and software requirements for Qwen 3.5

Before running Qwen 3.5 locally, you need to ensure your setup meets both hardware and software requirements for smooth inference. This guide will use an NVIDIA H200 GPU with 141GB of VRAM, combined with 240GB of system RAM, providing sufficient memory to efficiently run the MXFP4\_MOE version of Qwen 3.5 with MoE offloading enabled.

In general, for best performance, your total VRAM + RAM capacity should approximately equal the size of the quantization model you are downloading. Otherwise, llama.cpp might transfer to an SSD, but the inference process will be slower.

Regarding software, you need to install the latest NVIDIA GPU driver, along with a recent version of the CUDA Toolkit, to ensure full compatibility with llama.cpp and CUDA-accelerated inference.

## How to run Qwen 3.5 locally

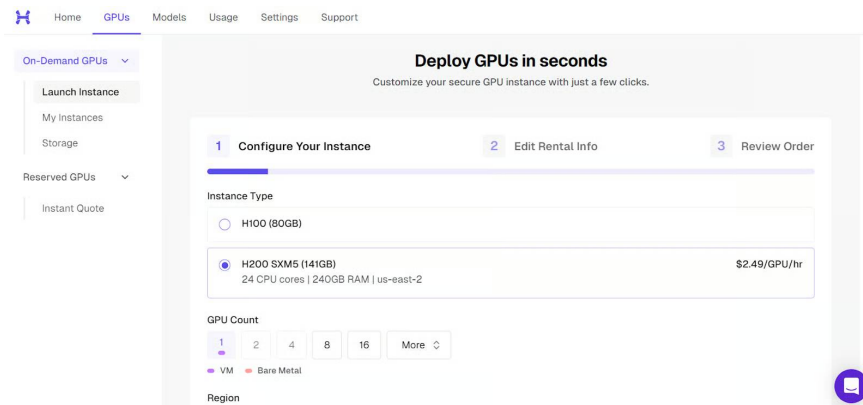
Now that you've met the prerequisites, let's look at a step-by-step guide on how to use Qwen 3.5 locally:

### 1. Set up the local environment.

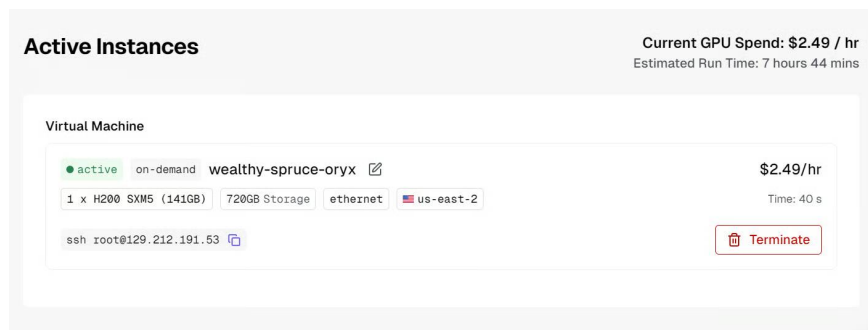
To run Qwen 3.5 locally, you need access to a computer with a powerful GPU. Since most laptops and desktops don't have enough VRAM or memory to handle models of this size, we'll be using a cloud-based GPU virtual machine.

This guide uses Hyperbolic to run the model privately. You can also use other providers such as RunPod, Vast.ai, or any GPU virtual machine platform you prefer. This article chooses Hyperbolic because it currently offers some of the most cost-effective GPU versions available.

Start by launching a new instance with a single H200 GPU.



After the machine boots up, you will see the public IP address and the SSH command needed to connect from your local terminal.

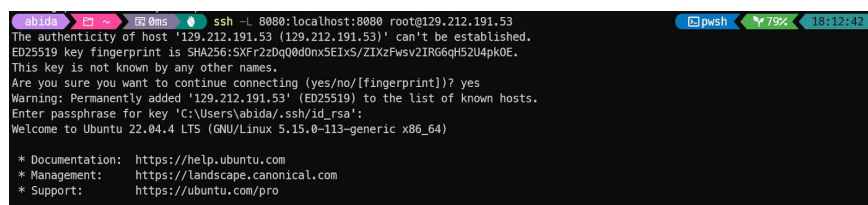


Before connecting, make sure you have set up local SSH and added your public SSH key when creating the virtual machine .

Once the instance is ready, connect to it via SSH with port forwarding. This is important because we want to access the local llama.cpp inference server through port 8080:

```
ssh -L 8080:localhost:8080 root@129.212.191.53
```

The first time you connect, type **yes** to confirm, then authenticate with your SSH key.



After logging in, verify that the GPU is correctly recognized:

```
nvidia-smi
```

You will see NVIDIA H200 listed in the results.

```
root@wealthy-spruce-oryx:~# nvidia-smi
Wed Feb 18 13:13:46 2026

+-----+
| NVIDIA-SMI 575.57.08              Driver Version: 575.57.08      CUDA Version: 12.9     |
+-----+-----+-----+-----+-----+-----+
| GPU  Name                   Persistence-M | Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp   Perf              Pwr:Usage/Cap |      Memory-Usage | GPU-Util  Compute M. |
|====+=====+====+=====+=====+=====+
|  0   NVIDIA H200                On          | 00000000:83:00.0 Off |             0         | |
| N/A   33C    P0                 78W / 700W   |  0MiB / 143771MiB |    0%      Default  |
|                                     |                 |                 | Disabled             |
+-----+-----+-----+-----+-----+-----+
|
| Processes:
| GPU   GI   CI          PID   Type   Process name                        GPU Memory
|  ID   ID   ID                   |                 |                               |     Usage
|-----+-----+-----+-----+-----+-----+
| No running processes found
|
+-----+-----+-----+-----+-----+-----+

```

Finally, install the necessary Linux packages to download, compile, and run llama.cpp:

```
sudo apt update sudo apt install pciutils build-essential cmake curl libcurl4-openssl-dev
```

Once complete, your environment is ready to install llama.cpp and run Qwen 3.5 locally.

## 2. Install llama.cpp with CUDA support.

llama.cpp is an open-source C and C++ inference engine that allows you to run large language models locally with minimal setup, supporting both CPU and GPU acceleration.

First, copy the llama.cpp archive:

```
git clone https://github.com/ggml-org/llama.cpp
```

Next, we configure the build to support CUDA with CMake. We enable CUDA using `-DGGML_CUDA=ON` and set the CUDA architecture to 90a since we're using NVIDIA H200 (Hopper class). This helps the build generate GPU code optimized for Hopper features.

```
cmake llama.cpp -B llama.cpp/build -DGGML_CUDA=ON -DCMAKE_BUILD_TYPE=Release
```

```

-- Detecting CUDA compiler ABI info - done
-- Check for working CUDA compiler: /usr/local/cuda/bin/nvcc - skipped
-- Detecting CUDA compile features
-- Detecting CUDA compile features - done
-- Using CMAKE_CUDA_ARCHITECTURES=90 CMAKE_CUDA_ARCHITECTURES_NATIVE=
-- CUDA host compiler is GNU 11.4.0
-- Including CUDA backend
-- gcm version: 0.9.7
-- gcm commit: ea003229d
-- Could NOT find OpenSSL, try to set the path to OpenSSL root folder in the system variable OPENSSL_ROOT_DIR (missing: OPENSSL_CRYPT
O_LIBRARY OPENSSL_INCLUDE_DIR)
CMake Warning at vendor/cpp-http/lib/CMakeLists.txt:150 (message):
  OpenSSL not found, HTTPS support disabled

-- Generating embedded license file for target: common
-- Configuring done
-- Generating done
-- Build files have been written to: /root/llama.cpp/build

```

Now compile the server binary file. llama-server is a built-in REST server that allows you to expose llama.cpp as an API endpoint:

```
cmake --build llama.cpp/build --config Release -j --clean-first --target llama-s
```

```

[ 95%] Building CXX object tools/server/CMakeFiles/server-context.dir/server-task.cpp.o
[ 95%] Building CXX object tools/server/CMakeFiles/server-context.dir/server-queue.cpp.o
[ 96%] Building CXX object tools/server/CMakeFiles/server-context.dir/server-common.cpp.o
[ 96%] Building CXX object tools/server/CMakeFiles/server-context.dir/server-context.cpp.o
[ 96%] Linking CXX static library libserver-context.a
[ 96%] Built target server-context
[ 96%] Generating loading.html.hpp
[ 98%] Generating index.html.gz.hpp
[ 98%] Building CXX object tools/server/CMakeFiles/llama-server.dir/server-http.cpp.o
[ 98%] Building CXX object tools/server/CMakeFiles/llama-server.dir/server.cpp.o
[ 98%] Building CXX object tools/server/CMakeFiles/llama-server.dir/server-models.cpp.o
[100%] Linking CXX executable ../bin/llama-server
[100%] Built target llama-server

```

Finally, copy the compiled binary files to the home directory for easy execution:

```
cp llama.cpp/build/bin/llama-* llama.cpp
```

### 3. Download the Qwen 3.5 model.

Now that you've installed llama.cpp, the next step is to download the actual Qwen 3.5 model weights in GGUF format. These files are quite large, so using the Hugging Face CLI is the most reliable way to load them directly onto your GPU.

Python needs to be installed first because Hugging Face's download tools and validation utilities are distributed as Python packages. Although llama.cpp itself is written in C++, Python makes managing model downloads and delivery much easier.

Start by installing pip:

```
sudo apt install python3-pip
```

Next, install the Hugging Face Hub client along with performance-enhancing tools. hf\_transfer and hf-xet significantly speed up downloads, which is crucial when downloading hundreds of gigabytes of model files.

```
pip -q install -U huggingface_hub hf-xet pip -q install -U hf_transfer
```

Now, let's download the Qwen 3.5 model from Hugging Face. In this tutorial, we'll only download the MXFP4\_MOE variant, which is optimized for efficient MoE inference:

```
hf download unsloth/Qwen3.5-397B-A17B-GGUF --local-dir models/Qwen3.5 --includ
```

```
root@wealthy-spruce-oryx:~# hf download unsloth/Qwen3.5-397B-A17B-GDUF \
--local-dir models/Qwen3.5 \
--include "MXFP4_MOE"
Downloading (incomplete total...): 3%█ | 6.81G/216G [00:06<01:50, 1.90GB/s]
Fetching 6 files: 17%█ | 1/6 [00:00<00:02, 1.67it/s]
```

Once the download is complete, the model files will be stored in `models/Qwen 3.5`, ready to be loaded into `llama.cpp` for local inference.

## 4. Launch the Qwen 3.5 model on a single GPU.

Now, we can launch Qwen 3.5 using `llama-server`. This provides us with an OpenAI-compatible endpoint API that can be called from local tools and applications.

Optimize the server for a single-GPU setup by doing three main things. First, enable the `--fit` option so that `llama.cpp` automatically balances the model between the GPU's VRAM and system RAM, instead of reporting an error when the model doesn't fit in all the VRAM .

Secondly, we use a larger context window with `--ctx-size 16384` so the server can handle longer prompts. Thirdly, we enable the `--jinja` option and pass `--chat-template-kwarg`s to control chat formatting and disable thought mode for faster and more direct responses.

Run the server with the command:

```
./llama.cpp/llama-server --model models/Qwen3.5/MXFP4_MOE/Qwen3.5-397B-A17B-MXFP4 --fit --ctx-size 16384 --jinja --chat-template-kwarg 'jinja'
```

While the model is loading, you'll see it using both GPU VRAM and system memory, which is normal for a large MoE model.

```
load_tensors: offloading output layer to GPU
load_tensors: offloading 59 repeating layers to GPU
load_tensors: offloaded 61/61 layers to GPU
load_tensors: CPU_Mapped model buffer size = 1030.62 MiB
load_tensors: CPU_Mapped model buffer size = 7862.31 MiB
load_tensors: CPU_Mapped model buffer size = 47639.83 MiB
load_tensors: CPU_Mapped model buffer size = 15792.86 MiB
load_tensors: CUDA0 model buffer size = 137627.10 MiB
.....|
```

Once the loading process is complete, the server will be accessible at:

1. `0.0.0.0:8080` on the virtual machine
2. `http://127.0.0.1:8080` on your local machine after SSH port forwarding

```
srv load_model: for more info see https://github.com/ggml-org/llama.cpp/pull/16391
init: chat template, example_format: '<|im_start|>system
You are a helpful assistant<|im_end|>
<|im_start|>user
Hello<|im_end|>
<|im_start|>assistant
Hi there<|im_end|>
<|im_start|>user
How are you?<|im_end|>
<|im_start|>assistant
<think>
'
srv init: init: chat template, thinking = 1
main: model loaded
main: server is listening on http://0.0.0.0:8080
main: starting the main loop...
srv update_slots: all slots are idle
```

Let the server continue running. On your local computer, open a new terminal window and reconnect using SSH port forwarding:

```
ssh -L 8080:localhost:8080 root@129.212.191.53
```

Next, test the server by listing the available models:

```
curl -s http://127.0.0.1:8080/v1/models
```

If you see Qwen 3.5 in the response, your server is running correctly and you are ready to call it from the OpenAI SDK and your local applications.

```
root@wealthy-spruce-oryx:~# curl -s http://127.0.0.1:8080/v1/models
{"models":[{"name":"Qwen3.5","model":"Qwen3.5","modified_at":"","size":"","digest":"","type":"model","description":"","tags":[],"capabilities":{"completion":true},"parameters":{"details":{"parent_model":"","format":"gguf","family":"","families":[""],"parameter_size":"","quantization_level":""}}],"object":"list","data":[{"id":"Qwen3.5","object":"model","created":1771421832,"owned_by":"llamacpp","meta":{"vocab_type":2,"n_vocab":248320,"n_ctx_train":262144,"n_embd":4096,"n_params":396346350336,"size":216143227904}}]}root@wealthy-s
```

## 5. Test the Qwen 3.5 model using the OpenAI SDK

Now that the Qwen 3.5 inference server is running, the next step is to verify that it works correctly with real-world client applications. One of the biggest advantages of llama.cpp is that llama-server provides an API compatible with OpenAI, meaning you can use the official OpenAI SDK without changing your code structure.

First, install the Python OpenAI package on your local computer (or inside a virtual machine if you prefer):

```
pip install openai
```

Now, let's run a simple test script. This script connects to your local forwarded endpoint at `http://127.0.0.1:8080/v1` instead of OpenAI's cloud server.

```
python3 - 'PY' from openai import OpenAI client = OpenAI( base_url="http://127.0
```

There are a few important details to understand here:

1. `base_url` points to your local Qwen 3.5 server, not the OpenAI API.
2. The API key is still required by the SDK, but llama.cpp doesn't require authentication, so any placeholder value will work.
3. The model name="Qwen 3.5" matches the alias set when the server starts up.

If everything is configured correctly, you will receive quick and clear feedback from the model.

```
root@wealthy-spruce-oryx:~/my-project# python3 - <<'PY'
from openai import OpenAI

client = OpenAI(
    base_url="http://127.0.0.1:8080/v1",
    api_key="sk-no-key-required"
)

response = client.chat.completions.create(
    model="Qwen3.5",
    messages=[
        {"role": "user", "content": "Write one sentence about AI agents."}
    ]
)

print(response.choices[0].message.content)
PY
AI agents are autonomous software entities that perceive their environment, reason about available information, and take actions to a
chieve specific goals without continuous human intervention.
```

This confirms that:

1. The Qwen 3.5 model has been successfully loaded.
2. The llama.cpp server is running correctly.
3. Your SSH port forwarding is working.
4. The endpoint is fully compatible with OpenAI-style applications.

At this point, you can integrate Qwen 3.5 into any local tool, agent workflow, or application that already supports the OpenAI API format.

## 6. Develop a text-based user interface (TUI) for stock trading using Llama.cpp's WebUI.

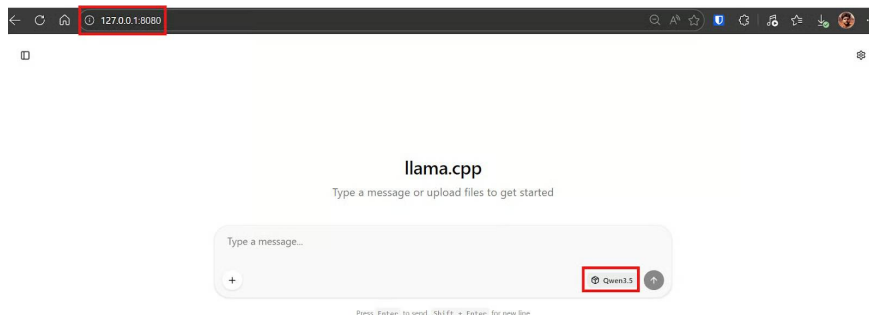
Llama.cpp includes a built-in WebUI, similar to ChatGPT, which you can use to chat directly with the model in your browser. This is useful for quick testing, iterative operations, and code generation without having to write any client scripts beforehand.

Because SSH port forwarding has been set up, you can open the WebUI on your local computer and it will function as if the server were running on your laptop.

By default, the WebUI is available at:

`http://127.0.0.1:8080`

If this page loads, it confirms two things: your SSH tunnel is working correctly, and the Qwen 3.5 server is locally accessible while still running privately on a GPU virtual machine.



After you enter the WebUI, paste this prompt. The goal is for the model to generate both Python code and a brief user guide.

Xây dựng một ứng dụng giao diện người dùng (TUI) để  
ng Python "Stock Screener Trainer" cho phép tôi nh  
?n rich (không phải giao diện web). Ứng dụng này cho phép tôi nh  
p danh sách mã chứng khoán, chọn lọc (theo trọng giá trị/cổ tức) và mức  
rủi ro (thấp/trung bình/cao), lấy các chỉ số công khai cho mỗi mã c  
? chứng khoán để phân tích, hiển thị trạng thái load theo từng loại, sau đó t  
?o một báo cáo và phân "Top 5 theo quy định của tôi" với lý do  
nh báo rõ ràng "chúng mang tính chất giáo dục, không phải là  
i khuyến tài chính", và lưu toàn bộ kết quả vào file 'results.csv'.

Within seconds, Qwen 3.5 will generate an `app.py` file, usually a brief explanation of how to run it.

```
BASH
pip install rich yfinance
```

2. Run the application:

```
BASH
python app.py
```

**Features Included**

1. **Interactive Input:** Prompts for tickers (comma-separated), Strategy Mode (Growth/Value/Dividend), and Risk Profile.
2. **Live Loading:** Uses `rich.progress` to show a spinner and progress bar while fetching data from Yahoo Finance.
3. **Scoring Logic:**
  - **Growth:** Rewards high EPS growth and higher Beta.
  - **Value:** Rewards low P/E, low P/B, and low Beta.
  - **Dividend:** Rewards high Dividend Yield and stability.

Type a message...

Qwen3.5

Now switch to your local terminal (laptop). Install the necessary libraries for the created application:

```
pip install rich yfinance
```

This operation installs:

1. Rich for text-based user interface (TUI) layouts, tables, prompts, and progress indicators.
2. Use yfinance to get free, publicly available stock indices.

Create a file named app.py, paste the code generated by the model, and run it:

```
python3 app.py
```

After running the script, you will see the text-based user interface (TUI) launch correctly in your terminal. The application will prompt you to enter the stock ticker symbol you want to analyze, along with your preferred screening mode and risk level.

For example, the author of the article conducted experiments with three popular stocks.

```
root@wealthy-spruce-oryx:~/my-project# python3 app.py
```

```
Stock Screener Trainer
```

```
Enter tickers (comma separated, e.g., AAPL, MSFT, TSLA): AAPL,MSFT,TSLA|
```

After a short loading period, the tool will return a complete table of stock market indicators, highlighting the results based on scoring rules and saving everything to the results.csv file.

🔥 TOP 5 PICKS

1. AAPL - Score: 25.0
2. MSFT - Score: 25.0
3. TSLA - Score: 25.0

Full Stock Analysis

Ticker	Name	Price	P/E	Div Yield %	Beta	Score
AAPL	Apple Inc.	\$265.55	33.66	39.00	1.11	25.0
MSFT	Microsoft Corporatio	\$401.24	25.09	92.00	1.08	25.0
TSLA	Tesla, Inc.	\$413.96	386.88	0.00	1.89	25.0

**DISCLAIMER:** This tool is for educational purposes only.  
This is NOT financial advice. Do your own due diligence before investing.

Saved results to results.csv

This is a great example of how Qwen 3.5 can create a fully functional application in a single go, using only a 4-bit quantized model endpoint and a simple prompt.

You finished reading the article "**How to run Qwen 3.5 locally on a single GPU**" edited by the [TipsMake](#) team. We hope this article has provided you with many useful tech tips and tricks. You can search for similar articles on tips and guides. Thank you for reading and for following us regularly.