# Ollama Cloud Models

26 October 2025      13:57

**Revise**

ollama → framework → LLM [locally].

↓

✔ Hardware. [Big LLM].

↓

2B, 5B, 7B

UM equipped with hardware
AWS, GCP, Runpods.
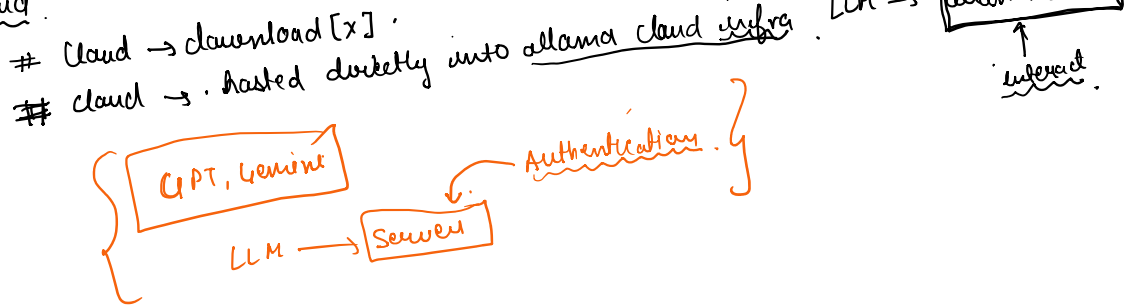
→ LLM → Server|Cloud
↓
RUN.

[ollama → Cloud → RUN very Big|Big models] → 6 7 B

✔ →

---

Ollama Cloud is a ==paid service== that allows ==you to run large, open-source language models on== ==powerful, datacenter-grade hardware== instead of on your own local computer.

Think of it as an extension of the ==local Ollama== you may be familiar with. ==While the standard Ollama is designed to run models on your personal Mac, Windows, or Linux machine,== ==Ollama Cloud== offloads the heavy computational work to remote servers equipped with powerful GPUs.

---

**ollama local**

→ Pull → menifest, Blobs, Config
→ Run → modelfile → LLM → W.M → Use.

---

**ollama Cloud**

\# Cloud → download [x].

\# cloud → hosted direktly into ollama cloud infra. LLM → [ollama cloud]
↑
interact.

{ GPT, Gemini
LLM → Server     Authentication }

---

⇒ ollama pull < ___-cloud >
→ Suffix → Cloud
→ Cloud → download [x] → [Cloud]
→ menifest, Config
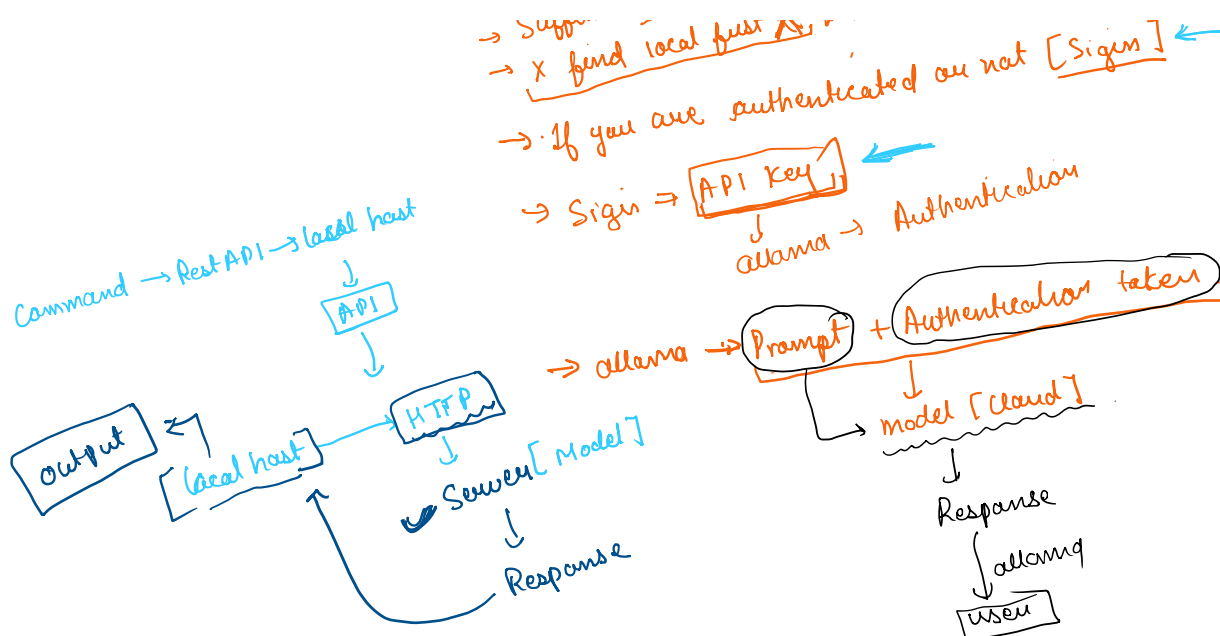→ Blob [Model weight] → Cloud Store.

---

⇒ ollama run < ___-cloud >     → http → [localhost! 11434]     → cloud
→ Suffix → Cloud → hosted on Cloud                 model =
→ x find local first ✗, ✗                          Prompt =
→ authenticated or not [Signin] ←

## Handwritten diagram (top)

Command → Rest API → local host
↓
API
↓
Output ← [local host] ← HTTP
↓
Server [Model]
↓
Response

→ Supp...
→ X find local first
→ If you are authenticated or not [Sign]
→ Sign → [API Key]
↓
allama → Authentication

→ allama → Prompt + Authentication token
↓
model [Cloud]
↓
Response
↓ allama
user

---

HTTP request → Rest API → [localhost: 11434]
↓
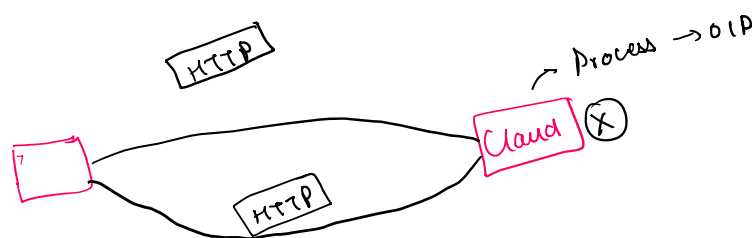Yes →

[Cloud] → local host

✓ [Cloud [X]] → [VM] → GPU

---

## What are the Benefits of Ollama Cloud?

The primary benefits of using Ollama Cloud revolve around **power, scale, and convenience**, especially when compared to running Ollama locally.

1. **Access to Extremely Large Models:** This is the main advantage. Many state-of-the-art open-source models (like those with hundreds of billions of parameters, such as deepseek-v3.1:671b) are simply too large to fit on the VRAM of consumer-grade GPUs. Ollama Cloud gives you access to these massive, powerful models without you needing to own the expensive hardware.

2. **Significantly Faster Performance:** Even if a model *can* run on your local machine, it might be very slow. Ollama Cloud uses high-performance datacenter GPUs, resulting in much faster inference speeds (quicker responses) and higher throughput (handling more requests at once).

3. **Seamless Transition from Local to Production:** You can develop and test an application on your laptop using a smaller, local model (like llama3) and then, when you're ready to deploy, switch to a more powerful cloud-hosted model (like gpt-oss:120b-cloud) with a one-line code change. The API and commands remain the same.

4. **No Hardware Management:** You don't need to buy, set up, or maintain expensive and power-hungry GPU hardware. Ollama handles all the infrastructure, so you can focus on building your application.

5. **Scalability for Applications:** A local laptop cannot reliably serve a production application with many concurrent users. Ollama Cloud is built to handle the scaling and availability required for real-world use cases.

6. **Privacy:** Similar to its local-first philosophy, Ollama states that it does not log or retain any user queries sent to the cloud service.

10,000

---

## Handwritten diagram (bottom)

HTTP
→ Process → O1P
[7] — [Cloud] (X)
HTTP

Gemini       Chatgp