

A Journey on Building AI Datacenter at Home

Affan Basalamah
IDNOG 10
2025

whoami

- Affan Basalamah
- @affanzbasalamah (X)
- affan@salamahsystems.com
- <https://s.id/salamahsystems.com>
- Consultant
- APNIC Community Trainer



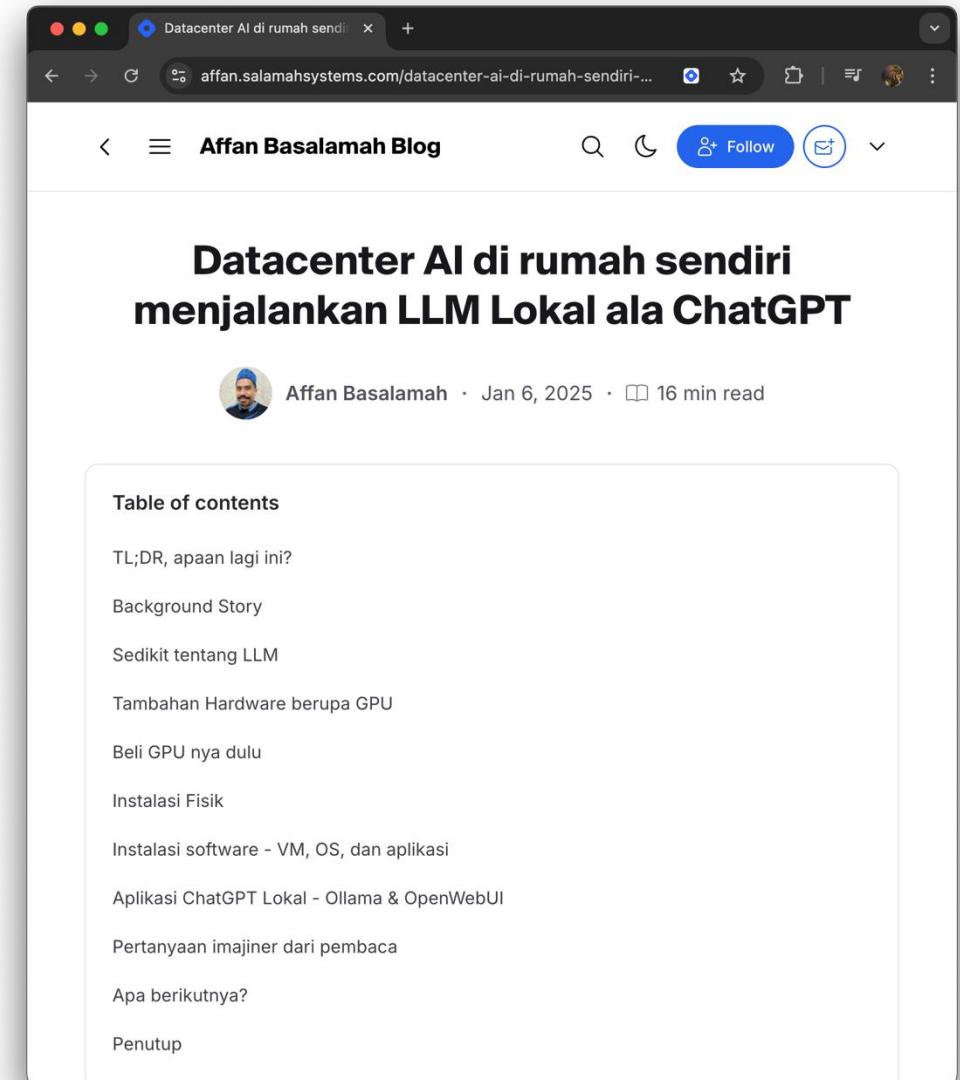
Story begins

- As a consultant, I want to have my own homelab to create PoC easily
 - PNETLAB, Containerlab
 - DNS, NMS, etc.
 - Local k8s, etc.
- And wrote a blog post about it:
<https://affan.salamahsystems.com/membuat-datacenter-di-rumah-sendiri>

The screenshot shows a web browser displaying a blog post. The title of the post is "Membuat Datacenter di Rumah Sendiri". The author is "Affan Basalamah", and the date is "Dec 14, 2024". The post has a "20 min read" duration. Below the title, there are three small icons: a heart, a brain, and a gear. The main content area is titled "Table of contents" and lists several sections: "Update: Bikin Datacenter AI dan DeepSeek lokal", "TL/DR, ceritanya apa ini?", "Background", "Survey, perencanaan, dan pembelian", "Instalasi dan Konfigurasi", "Penggunaan VM", and "Operational Expenditure (Opex) dan perjalanan nya". At the bottom of the content area, there is a "Show more" button. At the very bottom of the page, there is a footer bar with the text "Update: Bikin Datacenter AI dan DeepSeek lokal" and "IDNOG 10".

Rise of opensource AI and local AI homelab

- In 2024 every homelab youtubers create their own local AI
- Allow me to create one more case study for datacenter at home
- Experiment by installing GPU on server
- Tried the local AI chatbot, works
- Create another blogpost:
<https://affan.salamahsystems.com/datacenter-ai-di-rumah-sendiri-menjalankan-llm-lokal-ala-chatgpt>



**This is my journey building
AI Datacenter at Home**

Today's Discussion

- What is AI that we're talking about
- What can you do with AI
- How you build AI on your own datacenter at home

What is an AI?

AI as Large Language Model (LLM)

- Usually referred to as *Foundational Model*
 - Given the instruction, it can generate information (text) or multiple information (text, image, video: multimodal)



Source: O'Reilly Media – AI Engineering

Types of LLM in use

Model license:

- Closed model
 - Can only be accessed by API
- Open weight model
 - Can be downloaded only for weight
- Open source model
 - Can be downloaded all its data

LLM size & properties:

- Measured by billions (B) of parameters
- Base foundational model
 - Very large size (hundreds B of parameter, hundreds of GB size)
- Distilled model
 - Compressed the data from teacher model to student model
 - Between (1.5, 3, 7, 14, 32, 70 B of parameter, size between 1,5 to 80 GB size)
- Quantized model
 - Reduces the precision of model

LLM that you may know and use

Open weight model:

- DeepSeek-R1
- Google Gemma



Gemma 3

Open source model:

- Meta AI LLaMA 3.1
- Alibaba Qwen 3
- Mistral AI Mixtral-8x22B



Qwen3



MISTRAL
AI_

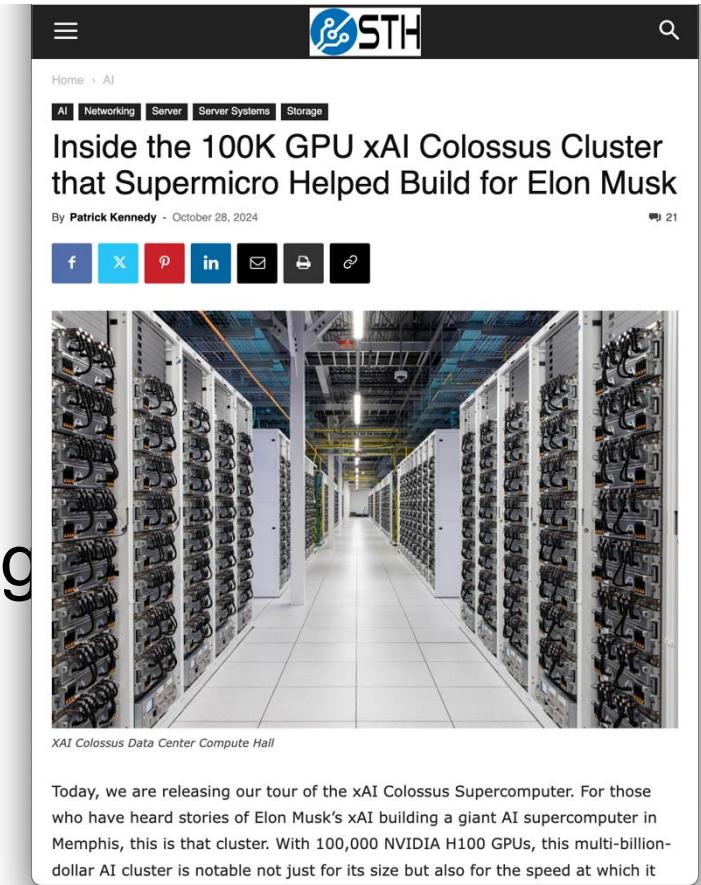
Closed model:

- OpenAI GPT
 - (4o, 4.1, o3, o4, etc.)
- Anthropic Claude
 - (Opus 4, Sonnet 4, Haiku 3.5)
- Google Gemini
 - (2.5 Flash, 2.5 Pro, etc.)
- xAI Grok



How to use the model

- **Training**
 - Build the foundational model from scratch
 - Training foundational model require very large
very large of data
- **Inferencing**
 - Use the trained model for AI application
 - We will use the trained model for AI datacenter at home
- **Finetuning**
 - Some model need to be finetuned for changing the characteristics of results



The screenshot shows a news article from STH (Supermicro Technology Hub) titled "Inside the 100K GPU xAI Colossus Cluster that Supermicro Helped Build for Elon Musk". The article is dated October 28, 2024, and includes social sharing icons for Facebook, X, LinkedIn, Email, and Print. Below the headline is a photograph of a massive server room filled with rows of server racks. A caption below the photo reads "XAI Colossus Data Center Compute Hall". The text of the article describes the cluster's size (100,000 NVIDIA H100 GPUs) and speed.

Inside the 100K GPU xAI Colossus Cluster that Supermicro Helped Build for Elon Musk

By Patrick Kennedy - October 28, 2024

f X p in e p

XAI Colossus Data Center Compute Hall

Today, we are releasing our tour of the xAI Colossus Supercomputer. For those who have heard stories of Elon Musk's xAI building a giant AI supercomputer in Memphis, this is that cluster. With 100,000 NVIDIA H100 GPUs, this multi-billion-dollar AI cluster is notable not just for its size but also for the speed at which it

How to host your own AI?

LLM needs to be run in systems with GPU (most of the time)

- LLM require heavy computation, basically matrix multiplication and tensor operation executed in parallel
- LLM weight in the system will be loaded to GPU memory to run
- Size of GPU memory (VRAM) will determine the size of the LLM
 - Bigger the model size, bigger the GPU VRAM needed
 - Finetuning require 2x-3x GPU VRAM compared with inferencing
 - Training require more, usually clustering the GPU systems
- There are more than that, but I try to simplify things
 - FLOPS, HBM size & speed, etc.

Where can I download LLM?

- You cannot download closed LLMs
- You can only download opensource/openweight LLMs
- <https://huggingface.co/> (HF) is the place to go
 - Place for complete opensource/openweight AI model
 - Model file usually with **.gguf** extension (*GPT-generated universal format*)
- For first timer, the fastest route is to use <https://ollama.com/>
 - Runtime to simplify LLM download in local system
 - Combine with OpenWeb-UI to create local ChatGPT
 - Can be used by AI apps (e.g. n8n, Dify) as local AI model

huggingface.co/deepseek-ai/DeepSeek-R1-0528

deepseek-ai/DeepSeek-R1-0528

like 2.31k Follow DeepSeek 83.1k

Text Generation Transformers Safetensors deepseek_v3 conversational custom_code

text-generation-inference fp8 arxiv:2501.12948 License: mit

Train Deploy Use this model

Model card Files Community 104

Edit model card

DeepSeek-R1-0528

deepseek

Downloads last month 377,015

Safetensors Model size 685B params Tensor type BF16 · F8_E4M3 · F32 Chat template Files info

Inference Providers NEW

Text Generation Examples

Input a message to start chatting with deepseek-ai/DeepSeek-R1-0528.

Your sentence here... Send

DeepSeek Homepage Chat DeepSeek R1

Hugging Face DeepSeek AI

Discord DeepSeek AI WeChat DeepSeek AI

X Twitter deepseek ai

License MIT

Paper Link

1. Introduction

The DeepSeek R1 model has undergone a minor version upgrade, with the current version being DeepSeek-R1-0528. In the latest update, DeepSeek R1 has significantly improved its depth of reasoning and inference capabilities by leveraging increased

huggingface.co/collections/deepseek-ai/DeepSeek-R1-678e1e131c01...

DeepSeek-R1 updated May 29

Upvote 756

deepseek-ai/DeepSeek-R1

Text Generation :: 685B Updated Mar 27 ↓ 824k ⚡ 12.5k

deepseek-ai/DeepSeek-R1-Zero

Text Generation :: 685B Updated Mar 27 ↓ 2.26k ⚡ 930

deepseek-ai/DeepSeek-R1-Distill-Llama-70B

Text Generation :: 71B Updated Feb 24 ↓ 144k ⚡ 706

deepseek-ai/DeepSeek-R1-Distill-Qwen-32B

Text Generation :: 33B Updated Feb 24 ↓ 411k ⚡ 1.42k

deepseek-ai/DeepSeek-R1-Distill-Qwen-14B

Text Generation :: 15B Updated Feb 24 ↓ 333k ⚡ 536

deepseek-ai/DeepSeek-R1-Distill-Llama-8B

Text Generation :: 8B Updated Feb 24 ↓ 999k ⚡ 775

deepseek-ai/DeepSeek-R1-Distill-Qwen-7B

Text Generation :: 8B Updated Feb 24 ↓ 483k ⚡ 678

deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B

Text Generation :: 2B Updated Feb 24 ↓ 965k ⚡ 1.27k

The Ollama homepage features a large 'DeepSeek-R1' logo at the top. Below it is a section titled 'Get up and running with large language models.' featuring a cartoon llama icon. A call-to-action button says 'Explore models →'. Below this, text mentions compatibility with macOS, Linux, and Windows. At the bottom, there's a navigation bar with links to Blog, Docs, GitHub, Discord, X (Twitter), Meetups, and Download.

© 2025 Ollama

The 'deepseek-r1' library page on ollama.com shows a list of available models. A red oval highlights the 'deepseek-r1:671b' row, which is labeled '404GB' in red text. The page includes a search bar with 'ollama run deepseek-r1' and a 'View all →' link.

Name	Size	Context	Input
deepseek-r1:latest	5.2GB	128K	Text
deepseek-r1:1.5b	1.1GB	128K	Text
deepseek-r1:7b	4.7GB	128K	Text
deepseek-r1:8b <small>latest</small>	5.2GB	128K	Text
deepseek-r1:14b	9.0GB	128K	Text
deepseek-r1:32b	20GB	128K	Text
deepseek-r1:70b	43GB	128K	Text
deepseek-r1:671b	404GB	160K	Text

Readme

IDNOG 10

Which LLMs require which GPU?

Model Size (full or distilled)	GPU type	Application Types
1.3B – 2B (e.g. DistilGPT2 , TinyLlama , Phi-2)	4 – 6GB GPU or CPU only	Embedded AI, lightweight chatbot
7B (e.g. LLaMA 2 7B , Mistral 7B , DeepSeek 7B , Gemma 7B)	8GB GPU	General-purpose assistants, RAG, code help
13– 14B (e.g. LLaMA 2 13B , StarCoder 14B)	24 – 32GB GPU	Smarter assistants with better reasoning, multi-turn chat, dev aid
30 – 34B (e.g LLaMA 2 30B , Mixtral , Yi-34B)	48– 80GB GPU	Complex assistants with solid reasoning, agents, domain-specific
65 – 70B (e.g. LLaMA 2 65B , Falcon 180B , Claude 2/3 class)	96 – 128GB GPU	GPT-4-class apps, enterprise copilots

NVid

NVIDIA Tesla K40

FP16: 5.046 TFLOPS Price-Performance Composite Score [?](#)

eBay Median Price - \$44.89 eBay Best Price - \$26.50

FP32: 5.046 TFLOPS
FP64: 1.682 TFLOPS
VRAM: 12 GB
Bandwidth: 288.4 GB/s

100.00% 

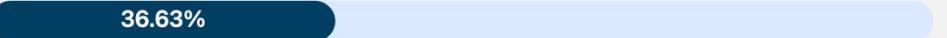
[Full Info](#)

NVIDIA Tesla M4

FP16: 2.195 TFLOPS Price-Performance Composite Score [?](#)

eBay Median Price - \$45.00 eBay Best Price - \$45.00

FP32: 2.195 TFLOPS
FP64: 0.069 TFLOPS
VRAM: 4 GB
Bandwidth: 88 GB/s

36.63% 

[Full Info](#)

NVIDIA Tesla M60

FP16: 4.825 x 2 TFLOPS Price-Performance Composite Score [?](#)

eBay Median Price - \$69.99 eBay Best Price - \$52.49

FP32: 4.825 x 2 TFLOPS
FP64: 0.151 x 2 TFLOPS
VRAM: 8 x 2 GB
Bandwidth: 160.4 x 2 GB/s

48.59% 

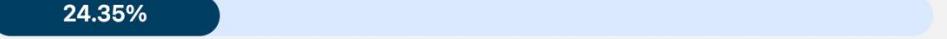
[Full Info](#)

NVIDIA Tesla P4

FP16: 0.089 TFLOPS Price-Performance Composite Score [?](#)

eBay Median Price - \$76.83 eBay Best Price - \$68.98

FP32: 5.704 TFLOPS
FP64: 0.178 TFLOPS
VRAM: 8 GB
Bandwidth: 192.3 GB/s

24.35% 

[Full Info](#)

<https://thedataaddi.com/hardware/gpucomp>

IDNOG 10

NVIDIA Tesla P4

Designed for AI inference and edge deployments, the Tesla P4 is compact and power-efficient, featuring 8GB of GDDR5 memory, and is ideal for recommendation systems, video processing, and low-latency inference tasks.

General Information	
Number Of GPUs	1
GPU Variant	GP104-895-A1
Architecture	Pascal
Process Size	16 nm
Release Date	2016-09-13
Generation	Tesla Pascal
Bus Interface	PCIe 3.0 x16
NVLink	Yes

Clock Speeds	
Base	1200 MHz

Performance Rates	
Pixel Rate	71.3 GPixel/s
Texture Rate	178.2 GTexel/s
FP16	0.089 TFLOPS
FP32	5.704 TFLOPS
FP64	0.178 TFLOPS

Memory Specifications	
VRAM Size	8 GB
VRAM Type	GDDR5
VRAM Bus	256 bit

GPU Prices Over Time

Price (\$)

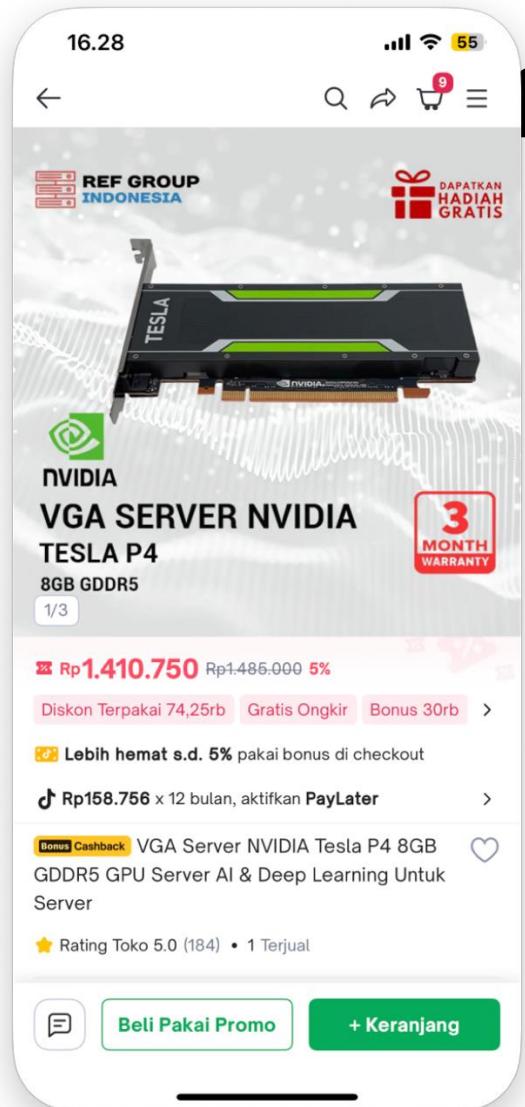
Best Price (Teal Line) | Median Price (Red Line)

Start Date: 13/11/2024 End Date: 20/07/2025

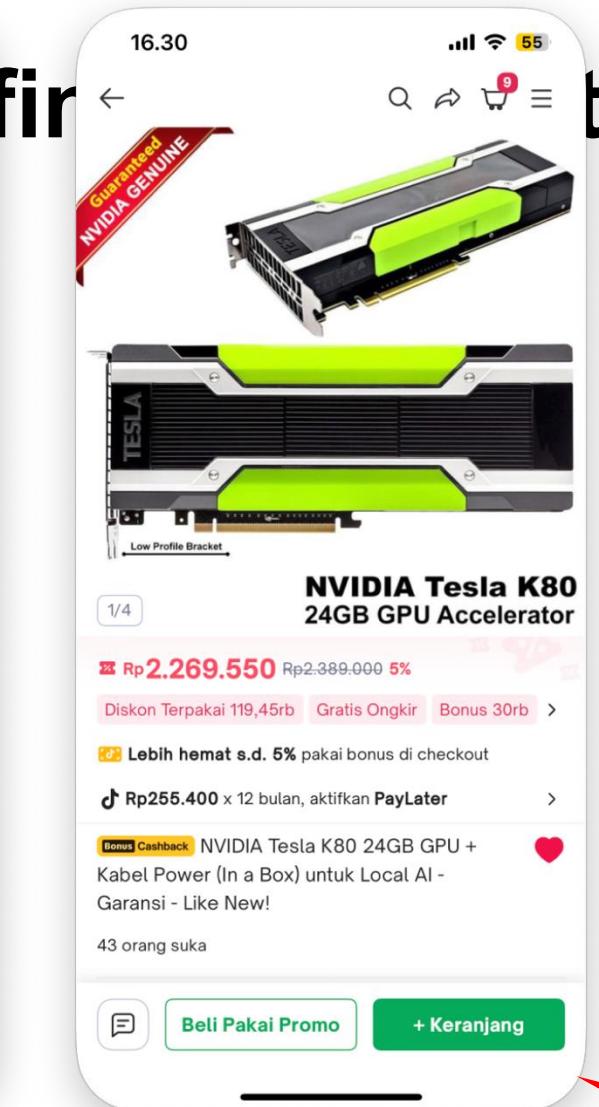
Apply Filter

Last Week Price History

Date Checked	Median Price	Best Price	Link
7/20/2025	\$76.83	\$68.98	ebay
7/19/2025	\$76.88	\$68.98	ebay
7/18/2025	\$76.88	\$68.98	ebay
7/17/2025	\$76.89	\$68.98	ebay



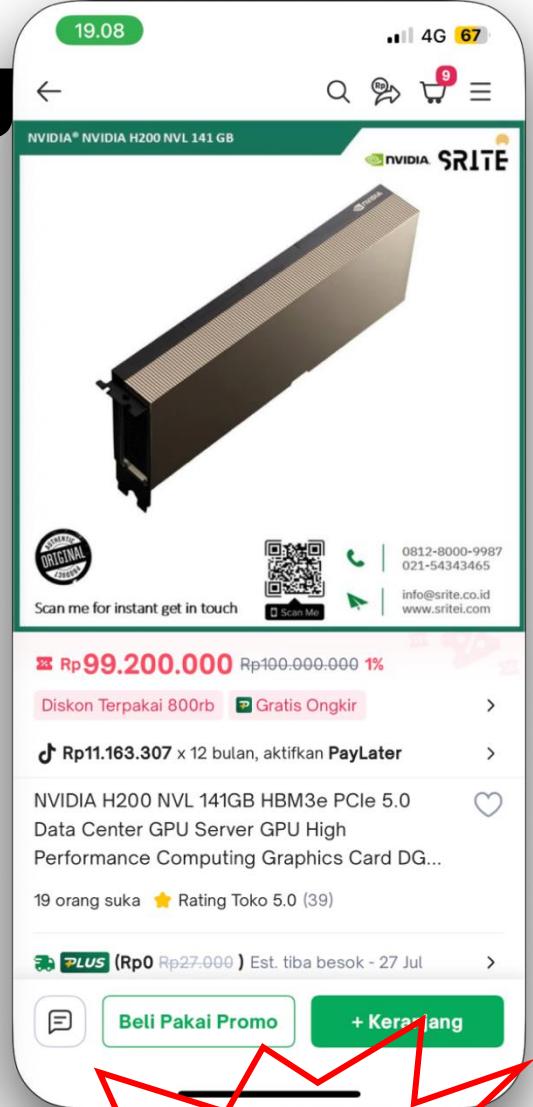
Tesla P4 – 8GB
Rp. 1.410.750



Tesla K80 – 2 x 12GB
Rp. 2.269.550



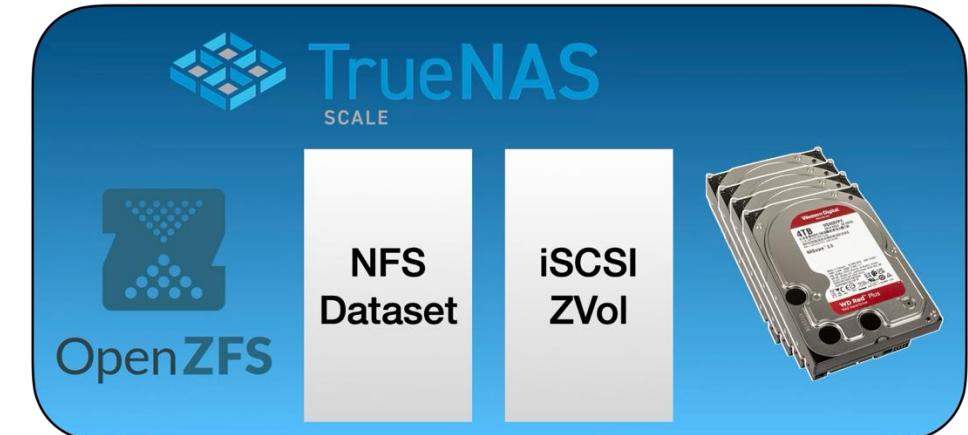
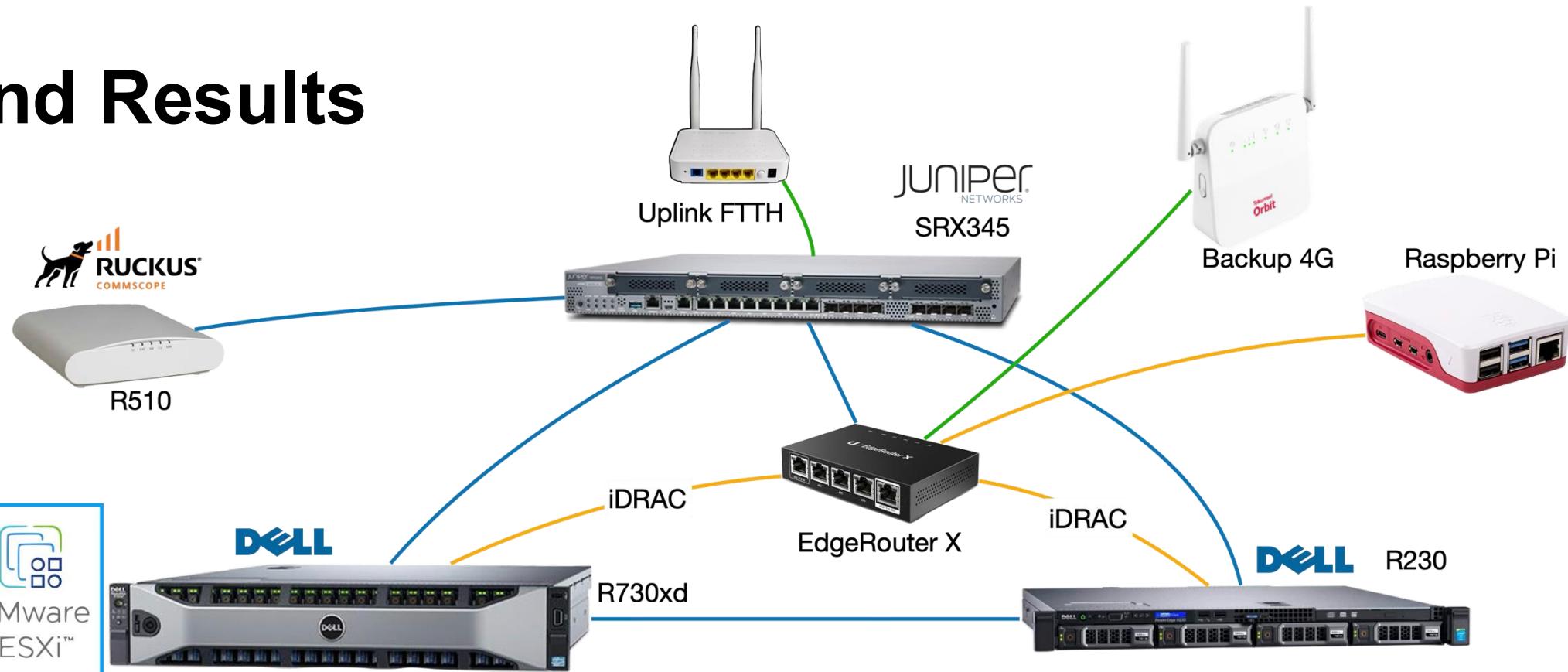
NVidia A100 – 80GB
Rp. 98.914.286



Nvidia H200 – 141GB
Rp. 99.200.000
IDMOS 10

How I do it in my home

End Results





Choose hardware for the systems

- For starter, better check aftermarket GPU
- Any PC or servers with a lot of RAM and PCIe slot
- Power depends on the GPU
 - Small form factor (SFF) GPU usually have lower power (75W), can be powered by PCIe slot itself
 - Dual width GPU usually have bigger power (300-500W), needs additional power cables

Tesla K80 – 24GB
Double width GPU



Tesla P4 – 8GB
SFF Single width GPU

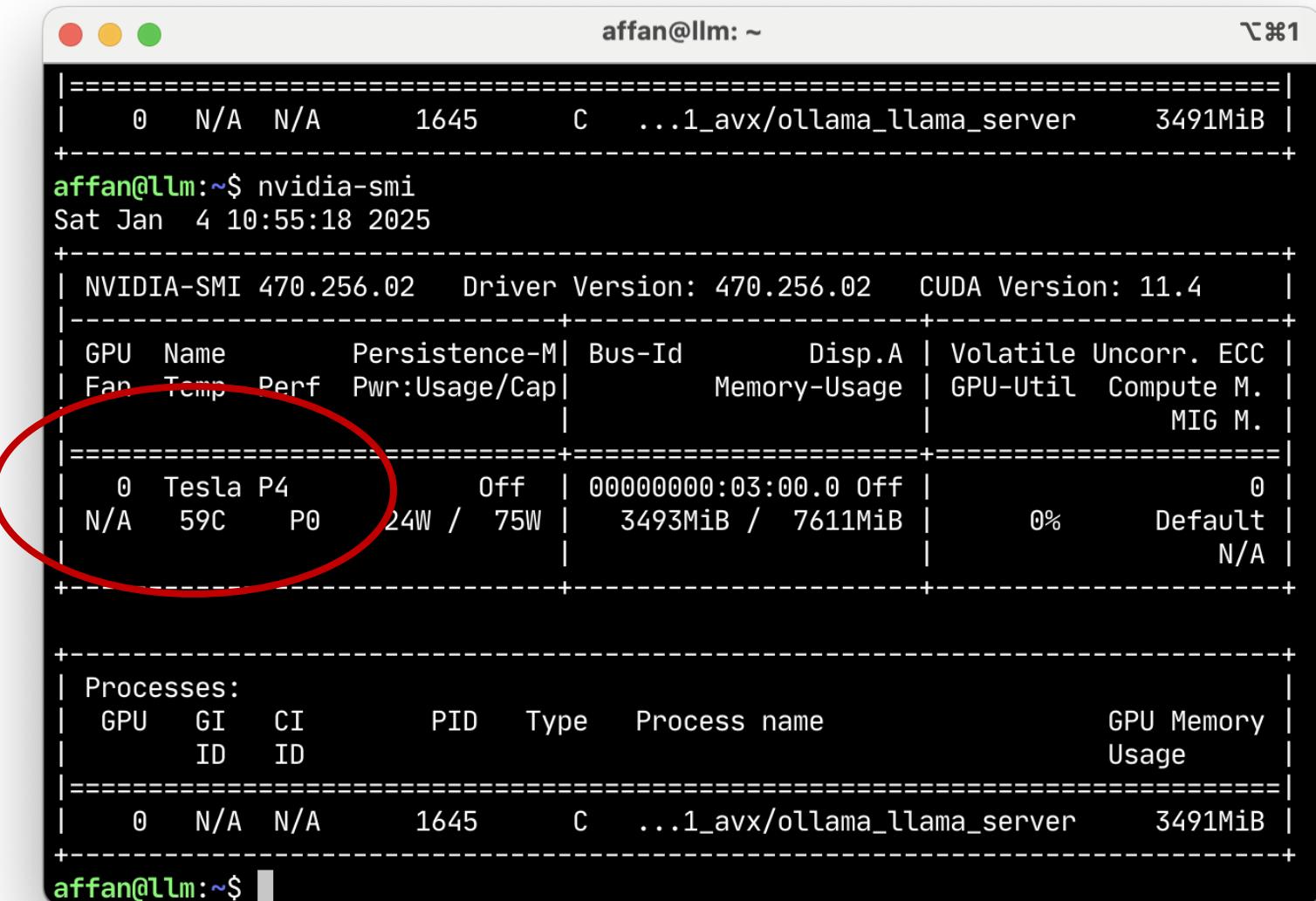


Software Stack to run AI Datacenter

- For desktop use use LMStudio <https://lmstudio.ai/>
 - Windows, macOS, Linux
- For server use Linux distro (myself use Ubuntu)
- Usually use Docker for easier installation
- GPU driver installed on OS
 - Learn how to install NVidia driver on Ubuntu
- Faster results: use integrated Ollama + OpenWeb-UI to create chatbot <https://openwebui.com/>
- Or better: use Ollama separately

Install NVidia Driver

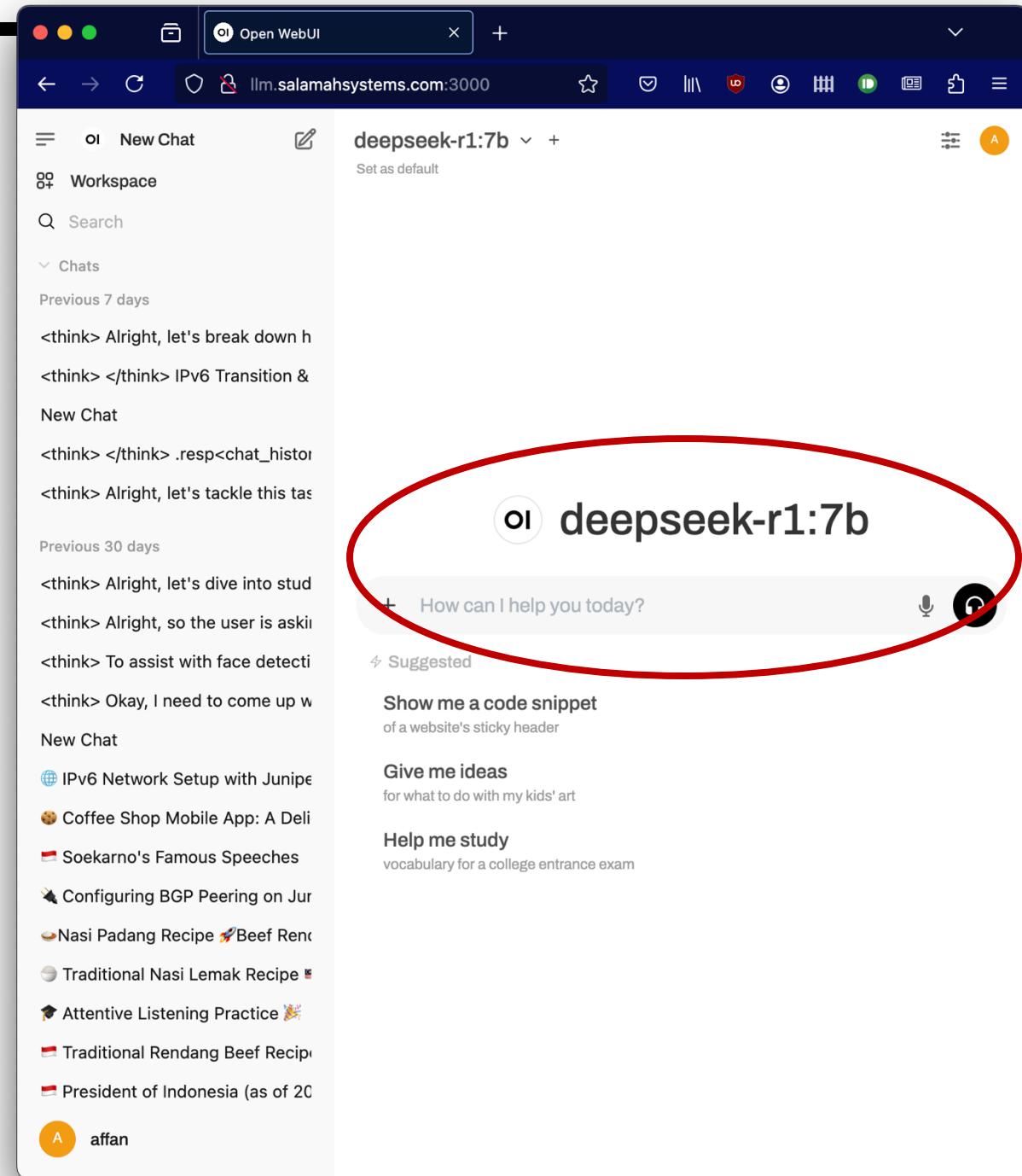
- Sometimes installation becomes a bit intricate
- “`sudo nvidia-smi`” and see your GPU model to begin install Ollama



```
affan@llm: ~
|=====
|   0  N/A N/A      1645      C ...1_avx/ollama_llama_server 3491MiB |
+---+
affan@llm:~$ nvidia-smi
Sat Jan  4 10:55:18 2025
+-----+
| NVIDIA-SMI 470.256.02    Driver Version: 470.256.02    CUDA Version: 11.4 |
+-----+
| GPU  Name     Persistence-M| Bus-Id     Disp.A  | Volatile Uncorr. ECC | | | | |
| Fan  Temp  Perf  Pwr:Usage/Cap| Memory-Usage | GPU-Util  Compute M.  |
|          |          |          |          |           |          |          MIG M. |
+-----+
| 0  Tesla P4          Off  | 00000000:03:00.0 Off |             0 | | | |
| N/A  59C   P0    24W / 75W | 3493MiB / 7611MiB | 0%     Default |
|          |          |          |          |           |          N/A |
+-----+
+-----+
| Processes:
| GPU  GI  CI          PID  Type  Process name          GPU Memory |
|          ID  ID
+-----+
| 0  N/A N/A        1645  C ...1_avx/ollama_llama_server 3491MiB |
+-----+
affan@llm:~$
```

Open WebUI

- Frontend for local LLM
- Suddenly you have local chatbot + RAG interface
 - Lots of function, check the docs
- Docker installation integrated with Ollama or separated

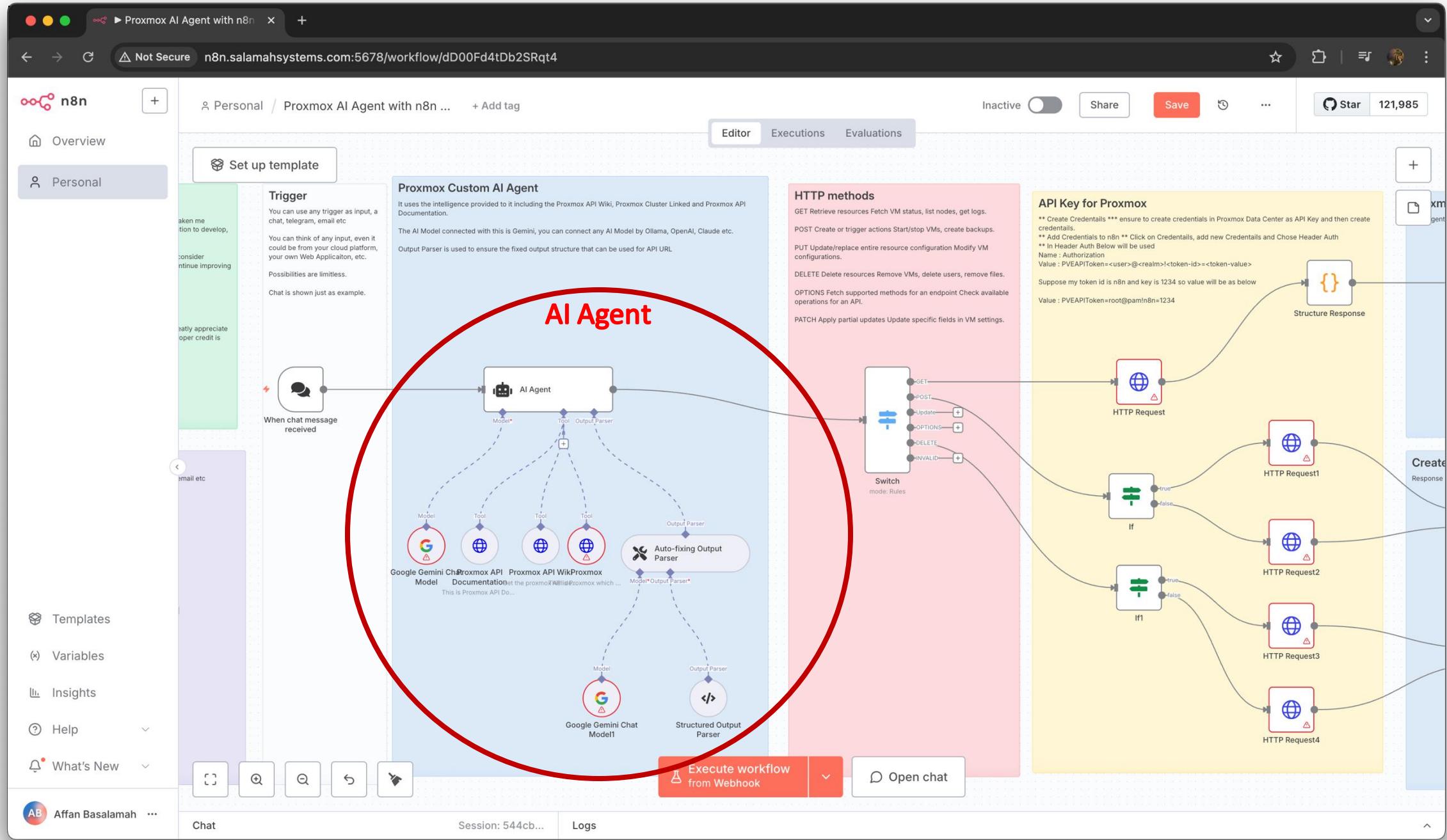


AI apps running on local LLM installation

- Local AI means private or secret document stay home and don't go to the remote LLM provider

AI Apps:

- Chatbot (ChatGPT local clone) with OpenWeb-UI
- Retrieve knowledge from documents using RAG
- Workflow automation from n8n accessing local Ollama instances



Chat with local LLMs using n8n

n8n.salamahsystems.com:5678/workflow/MVrC8a67s9XiedbD

Inactive Share Saved ...

Star 121,987

Overview Personal Chat with local LLMs using ... + Add tag

Editor Executions Evaluations

+

Chat with local LLMs using n8n and Ollama

This n8n workflow allows you to seamlessly interact with your self-hosted Large Language Models (LLMs) through a user-friendly chat interface. By connecting to Ollama, a powerful tool for managing local LLMs, you can send prompts and receive AI-generated responses directly within n8n.

How it works

1. When chat message received: Captures the user's input from the chat interface.
2. Chat LLM Chain: Sends the input to the Ollama server and receives the AI-generated response.
3. Delivers the LLM's response back to the chat interface.

Set up steps

- Make sure Ollama is installed and running on your machine before executing this workflow.
- Edit the Ollama address if different from the default.

Connect to remote Ollama

When chat message received

Chat LLM Chain

Model*

Ollama Chat Model

Ollama setup

- Connect to your local Ollama, usually on <http://localhost:11434>
- If running in Docker, make sure that the n8n container has access to the host's network in order to connect to Ollama. You can do this by passing `--net=host` option when starting the n8n Docker container

Open chat

AB Affan Basalamah ...

Chat Session: f6ff7... Logs

```
graph LR; Start(( )) --> ChatInput[When chat message received]; ChatInput --> ChatLLM[Chat LLM Chain]; ChatLLM --> OllamaModel((Ollama Chat Model)); OllamaModel --> End(( ));
```

The Journey Continues

Countless potential opportunities in local AI

- Lots of promises of network automation can be done with local AI help
- Automation tools (n8n) doesn't yet have connector for popular network OS/appliances
- Opportunities for opensource centralized SoT (source of truth) (e.g. Netbox, Nautobot)
- Lots more...

Do you want to join me in this journey?

- Join me at <https://s.id/salamahsystems.com>
 - Fill me a little survey
 - Maybe join a group
 - Maybe you want to receive newsletter

Thanks!