Generative AI with Google Clad

7 Days ⟶ ① Generative AI with Google

② Vertex AI ⟶ [Generative AI]

③ Build transformative LLM powered Application on google cloud.

④ Gemini on vertex AI

⑤ How to Implement Application with Gemini

⑥ RAG Implementation with Gemini pro, Vector search

⑦ Fine Tuning Foundation model on GCP

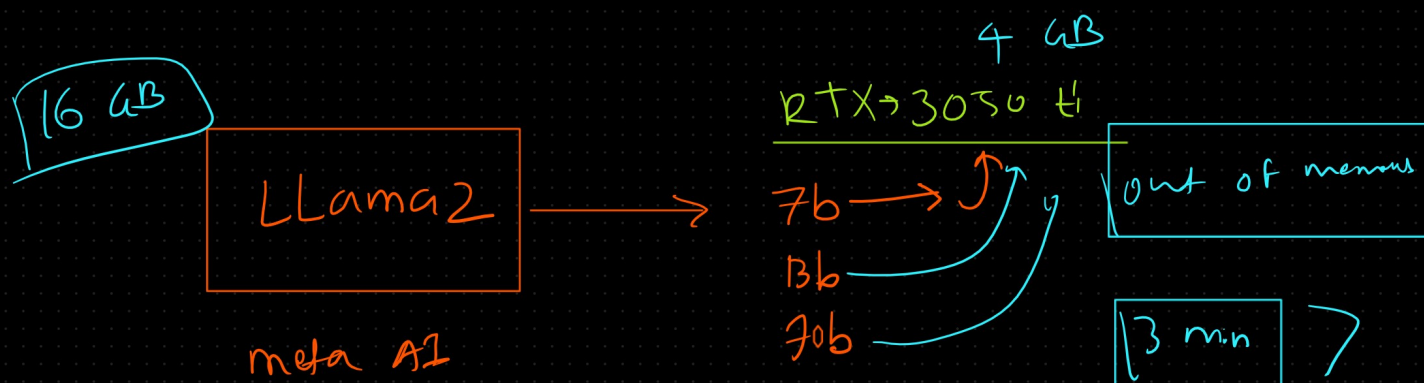# prerequisite:

→ python
→ Basics of Generative AI

   ├ LLM
   └→ Langchain
   └→ Vector DB
   └→ GCP Account — 300$

→ Langchain
→ llama Index
→ open source LLm ⟶⟶ 7 Gpu Based machine

4 GB

16 GB

RTX → 3050 ti

LLama2 ⟶ 7b ⟶ 😊 y | out of memory

meta AI

13b
70b

3 min 7

LLmops ⟶ Deploy open LLm

AWS ⟶ Sazemeketi
DLC

API ⟵ API Gateway
Lamda

openAI ⟶ GPT-3,4 → API key Req⇄ m
Res

| Model Name | Model ID | Max Total Tokens | Default Instance Type |
|---|---|---|---|
| Llama-2-7b | meta-textgeneration-llama-2-7b | 4096 | ml.g5.2xlarge |
| Llama-2-7b-chat | meta-textgeneration-llama-2-7b-f | 4096 | ml.g5.2xlarge |
| Llama-2-13b | meta-textgeneration-llama-2-13b | 4096 | ml.g5.12xlarge |
| Llama-2-13b-chat | meta-textgeneration-llama-2-13b-f | 4096 | ml.g5.12xlarge |
| Llama-2-70b | meta-textgeneration-llama-2-70b | 4096 | ml.g5.48xlarge |
| Llama-2-70b-chat | meta-textgeneration-llama-2-70b-f | 4096 | ml.g5.48xlarge |

## Google

2017 ⟶ Transformers

2018 ⟶ BERT

2019 → T5 ⟶ Text to Text Transfer Transformer
um 10B p model

2020 → LaMDA ⟶ LLM

2021 ⟶ AlphaFold

20222 → Palm —

2023 → Bend —

2024 →

## Generative AI

$\downarrow$

mcp

$\downarrow$

Vertex AI

Generative AI