

Crash Course: Open-Source LLMs



Mikiko Bazeley,
Head of Developer Relations
@Labelbox

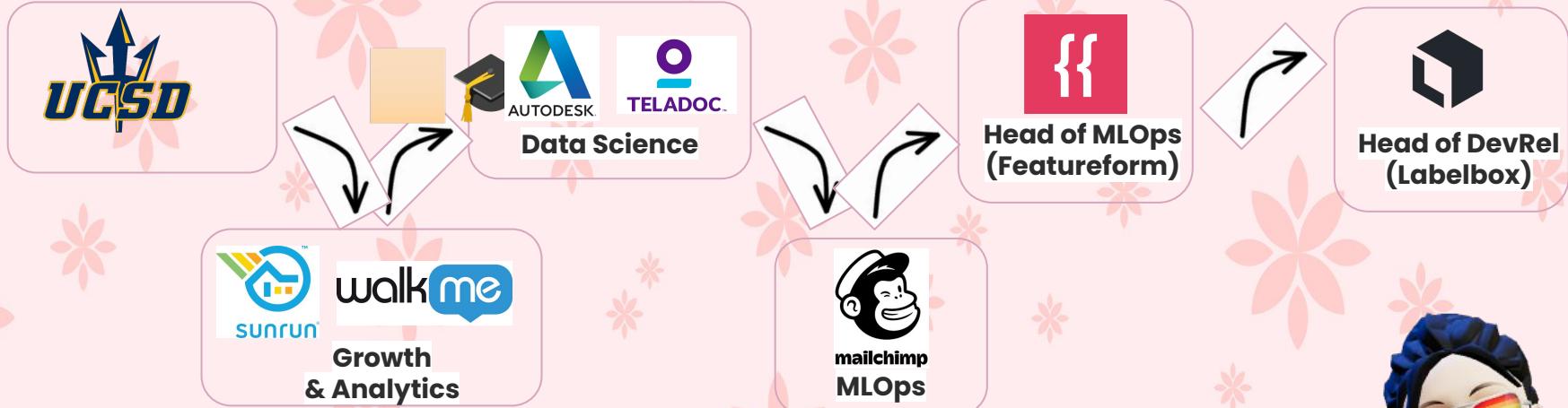


Objectives

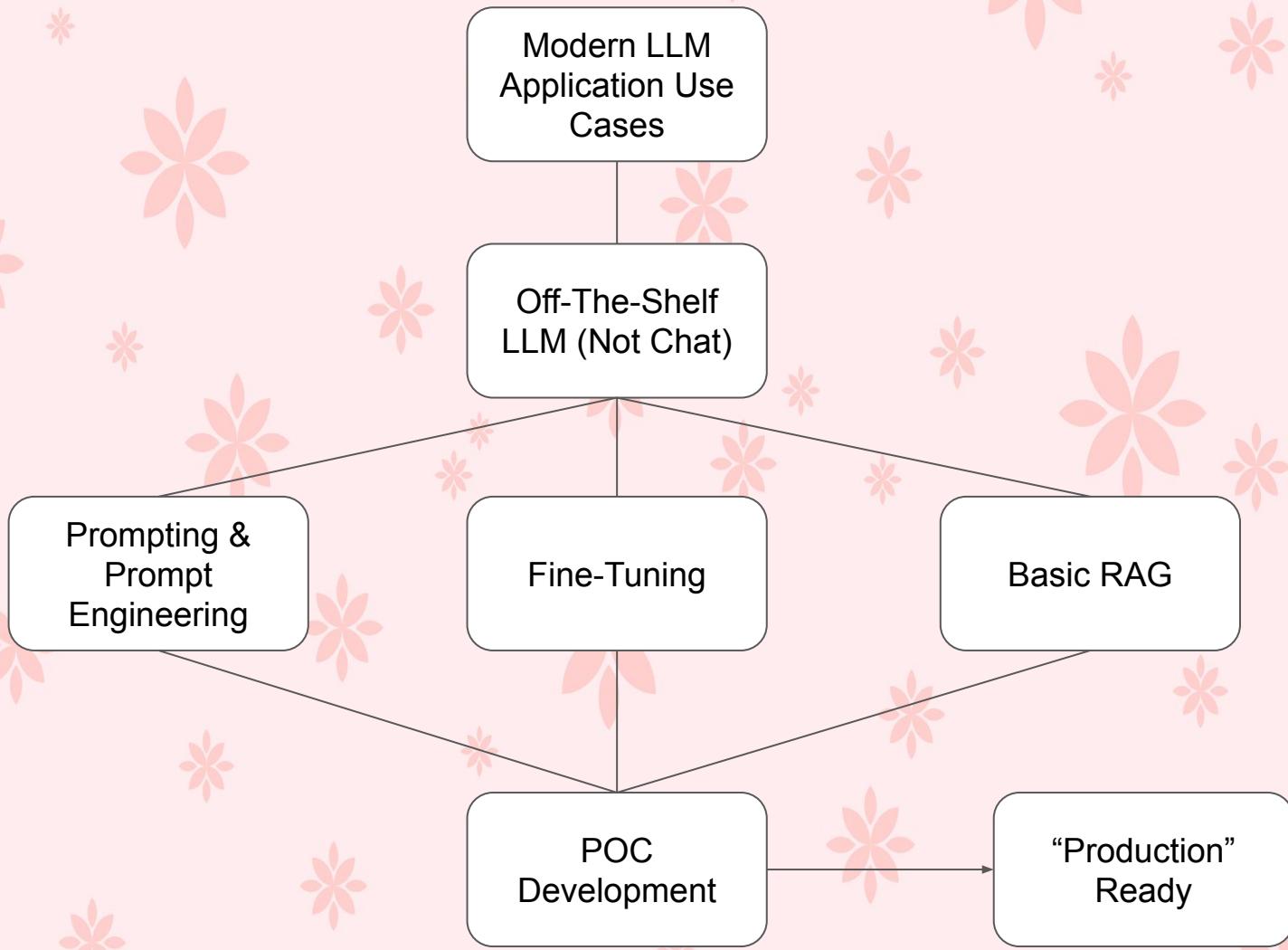
What we'll cover today:

- Main ways LLMs (both open-source & proprietary) are being used in applications
- Differences between OSS & closed-source models
- Introduce & define some key players & tools in the ecosystem
- Try to run some code samples (if not, we'll at least talk through them so you can run them outside the session)
- Answer Questions: What do you want to know about build LLM applications?

My Background



How We'll Proceed



How We'll Proceed

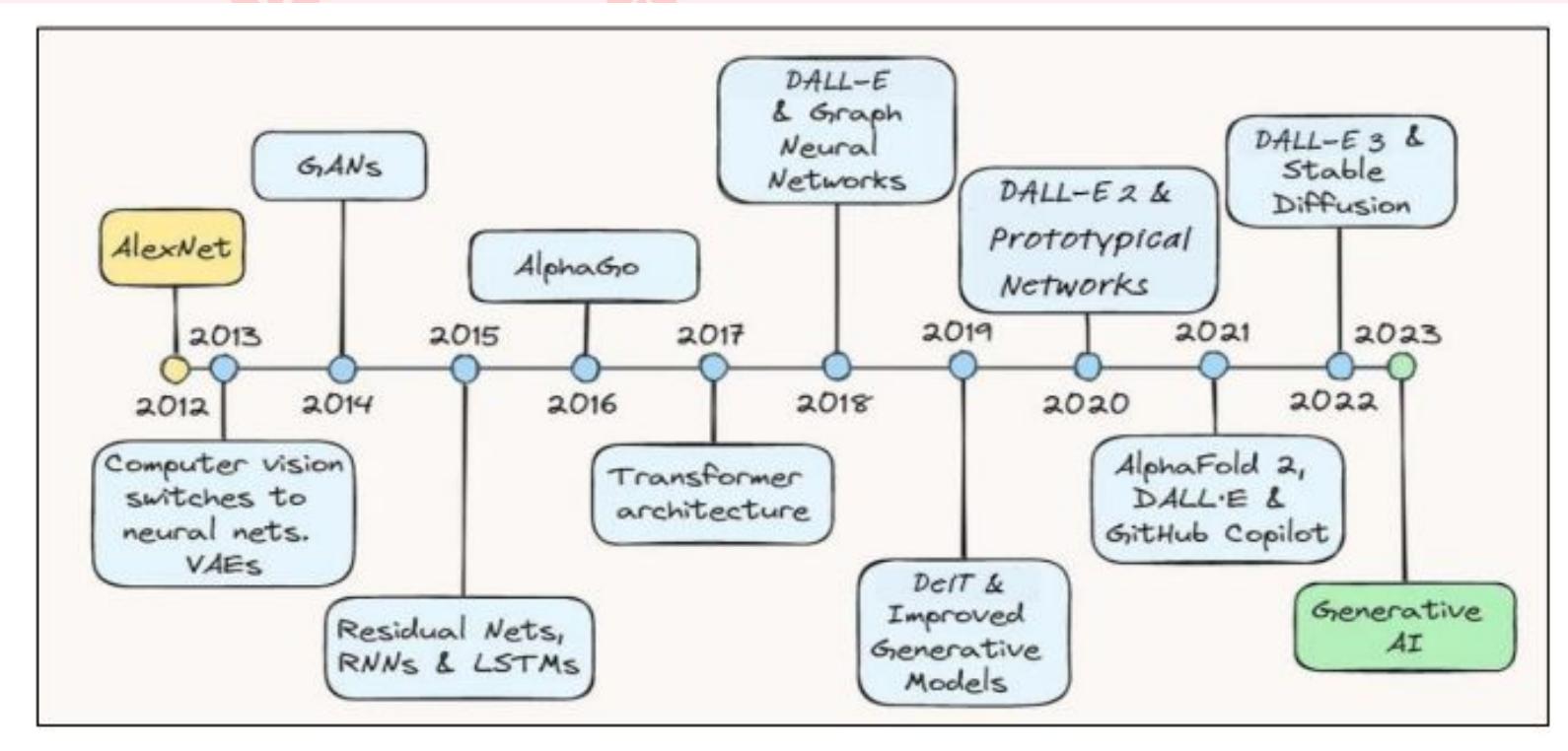
👉 Resources for all talks will be at this link 👈

https://github.com/MMBazel/LO_GenAI_Workshops

👉 Resources for OSS Models will be under /oss-models 👈

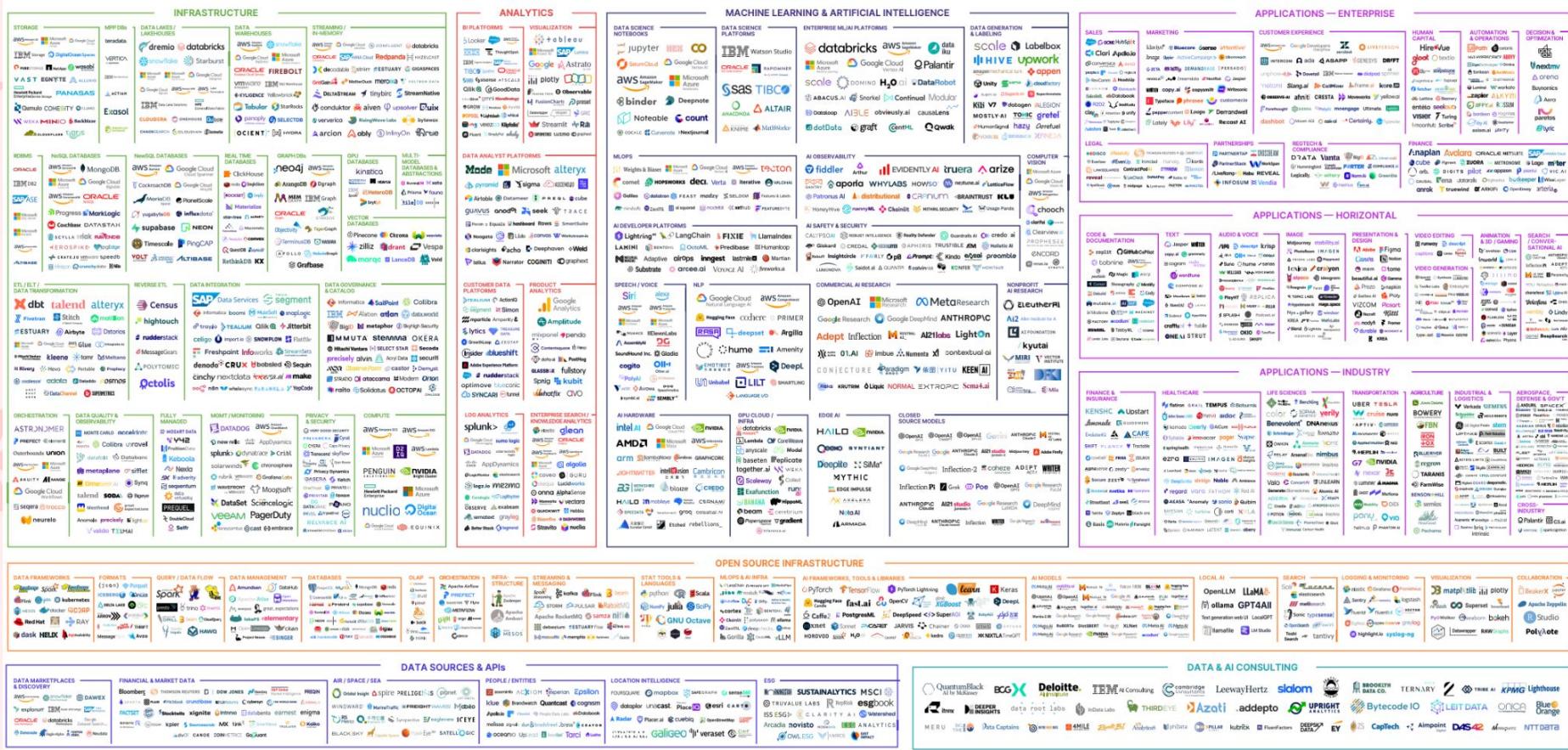
https://github.com/MMBazel/LO_GenAI_Workshops/tree/main/oss-models

We've Come A Long Way



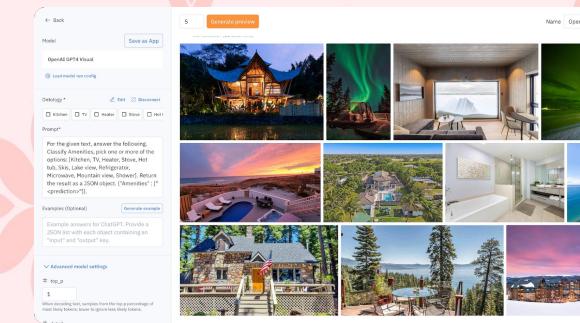
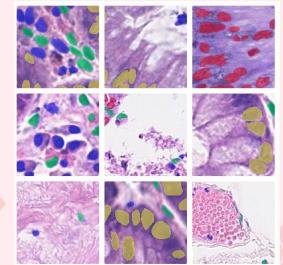
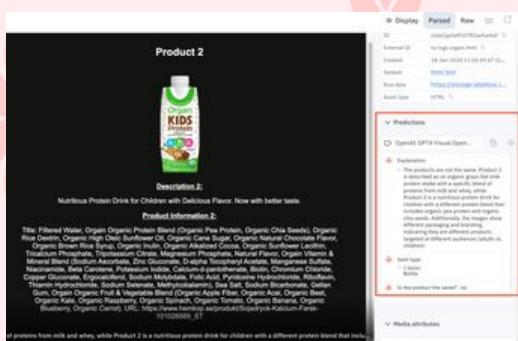
Almost Too Far???

THE 2024 MAD (MACHINE LEARNING, ARTIFICIAL INTELLIGENCE & DATA) LANDSCAPE



- Note: Diagram doesn't matter as much as the analysis does <https://mattturck.com/mad2024/>

Language as an interface is in every Gen-AI App



And is used in multiple modalities

Input/Output types:

-  : Text
-  : Image
- `</>` : Code
-  : Software tool use (text, code generation & execution)
-  : Video
-  : Music
-  : 3D
-  : Robot state

Model types:

-  →  : LLMs
-  +  →  : Multimodal LLMs
-  +  +  →  : Multimodal LLMs for Robotics
-  → `</>` : Text to Code
-  →  : Text to Software tool use
-  →  : Text to Image
-  →  : Text to Video
-  →  : Text to Music
-  →  : Image to 3D
-  →  : Text to 3D

There are currently more than 18.9K Text-to-Image generative models listed on Hugging Face

The screenshot shows the Hugging Face homepage with a search bar at the top. Below it, there's a navigation bar with links for Tasks, Libraries, Datasets, Languages, Licenses, and Other. The main content area is titled "Models" and shows a list of 447 models. Each model entry includes the name, a small icon, a brief description, and some statistics like file size and number of stars. The models listed are from various organizations like Salesforce, Microsoft, and others, and include names such as "Salesforce/blip-image-captioning-large", "llava-hf/llava-1.5-b-hf", and "microsoft/trocr-base-handwritten".

Table 1: Summary of Video Generation.

Model name	Year	Backbone	Task	Group
Imagen Video[29]	2022	Diffusion	Generation	Google
Pix2Seq-D[160]	2022	Diffusion	Segmentation	Google Deepmind
FDM[161]	2022	Diffusion	Prediction	UBC
MaskViT[162]	2022	Masked Vision Models	Prediction	Stanford, Salesforce
CogVideo[163]	2022	Auto-regressive	Generation	THU
Make-a-video[164]	2022	Diffusion	Generation	Meta
MagicVideo[165]	2022	Diffusion	Generation	ByteDance
TATS[166]	2022	Auto-regressive	Generation	University of Maryland, Meta
Phenaki[167]	2022	Masked Vision Models	Generation	Google Brain
Gen-I[168]	2023	Diffusion	Generation, Editing	RunwayML
LFDM[140]	2023	Diffusion	Generation	PSU, UCSD
Text2video-Zero[169]	2023	Diffusion	Generation	Picsart
Video Fusion[170]	2023	Diffusion	Generation	USAC, Alibaba
PyCo[34]	2023	Diffusion	Generation	Nvidia
Video LDM[36]	2023	Diffusion	Generation	University of Maryland, Nvidia
RIN[171]	2023	Diffusion	Generation	Google Brain
LVD[172]	2023	Diffusion	Generation	UCB
Dreamix[173]	2023	Diffusion	Editing	Google
MagicEdit[174]	2023	Diffusion	Editing	ByteDance
Control-A-Video[175]	2023	Diffusion	Editing	Sun Yat-Sen University
StableVideo[176]	2023	Diffusion	Editing	ZJU, MSRA
Tune-A-Video[78]	2023	Diffusion	Editing	NUS
Rerender-A-Video[177]	2023	Diffusion	Editing	NTU
Pix2Video[178]	2023	Diffusion	Editing	Adobe, UCL
InstructVid2Vid[179]	2023	Diffusion	Editing	ZJU
DiffAct[180]	2023	Diffusion	Action Detection	University of Sydney
DiffPose[181]	2023	Diffusion	Pose Estimation	Jilin University
MAGVIT[182]	2023	Masked Vision Models	Generation	Google
AnimateDiff[138]	2023	Diffusion	Generation	CUHK
MAGVIT V2[47]	2023	Masked Vision Models	Generation	Google
Generative Dynamics[183]	2023	Diffusion	Generation	Google
VideoCrafter[81]	2023	Diffusion	Generation	Tencent
Zeroscope[184]	2023	-	Generation	EasyWithAI
ModelScope	2023	-	Generation	Damo
Gen-2[23]	2023	-	Generation	RunwayML
Pika[22]	2023	-	Generation	Pika Labs
Emu Video[185]	2023	Diffusion	Generation	Meta
PixelDance[186]	2023	Diffusion	Generation	ByteDance
Stable Video Diffusion[27]	2023	Diffusion	Generation	Stability AI
W.A.L.T[187]	2023	Diffusion	Generation	Stanford, Google
Fairy[188]	2023	Diffusion	Generation, Editing	Meta
VideoPoet[189]	2023	Auto-regressive	Generation, Editing	Google
LGVI[190]	2024	Diffusion	Editing	PKU, NTU
Lumiere[191]	2024	Diffusion	Generation	Google
Sora[3]	2024	Diffusion	Generation, Editing	OpenAI

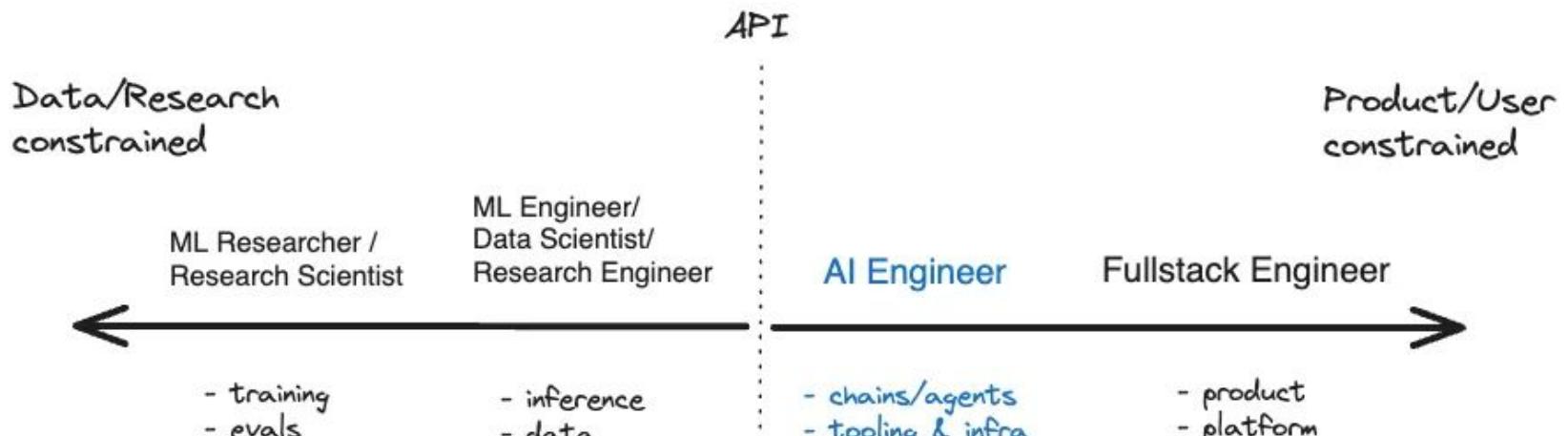
In order to build a career as an AI Engineer, you need to understand the low-level “why”

Someone else's POC
Gen-AI App
on Twitter

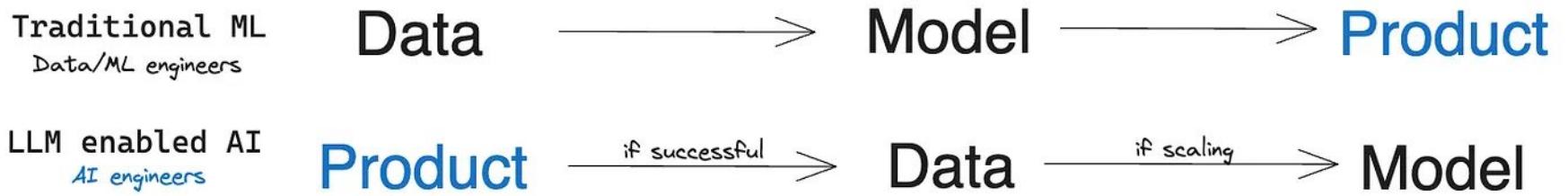


Your results

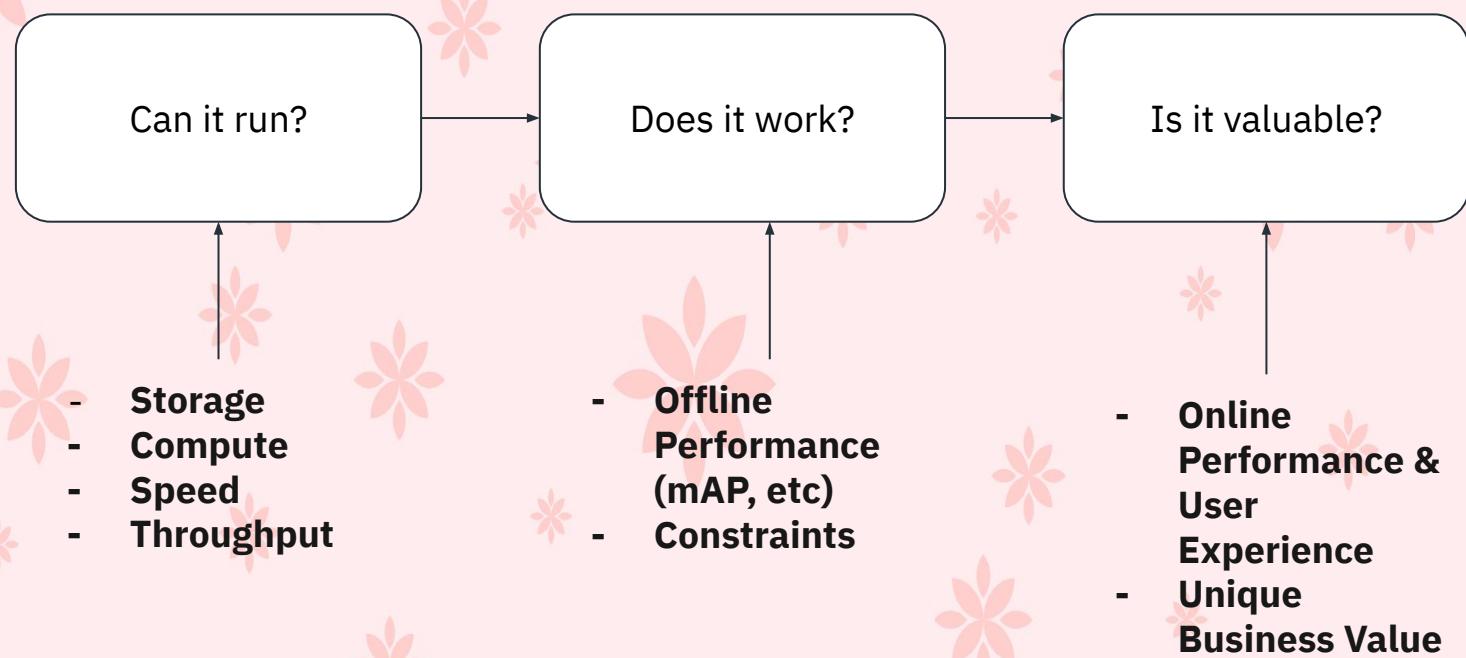
In order to build a career as an AI Engineer, you need to understand the low-level “why”



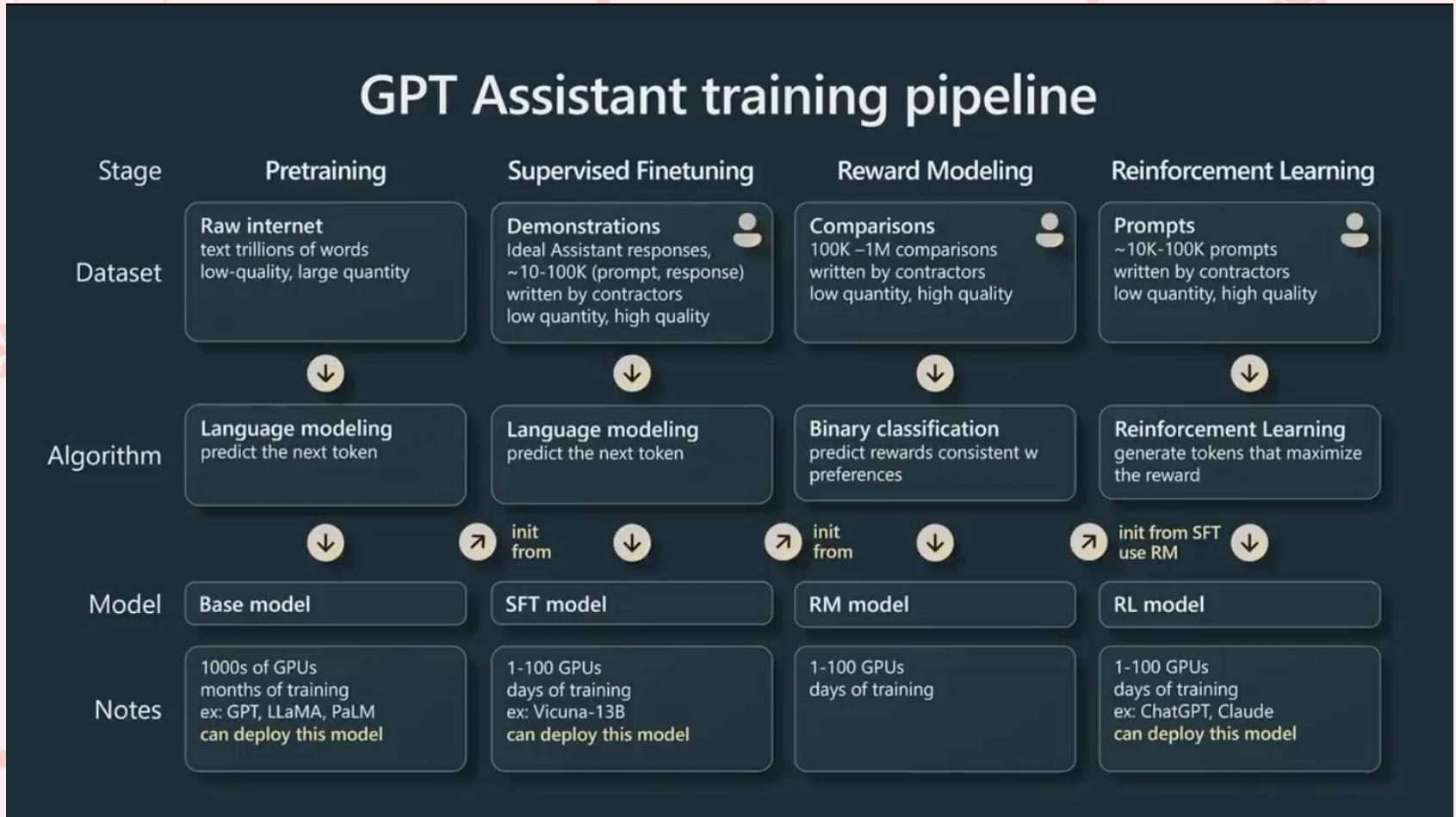
In order to build a career as an AI Engineer, you need to understand the low-level “why”



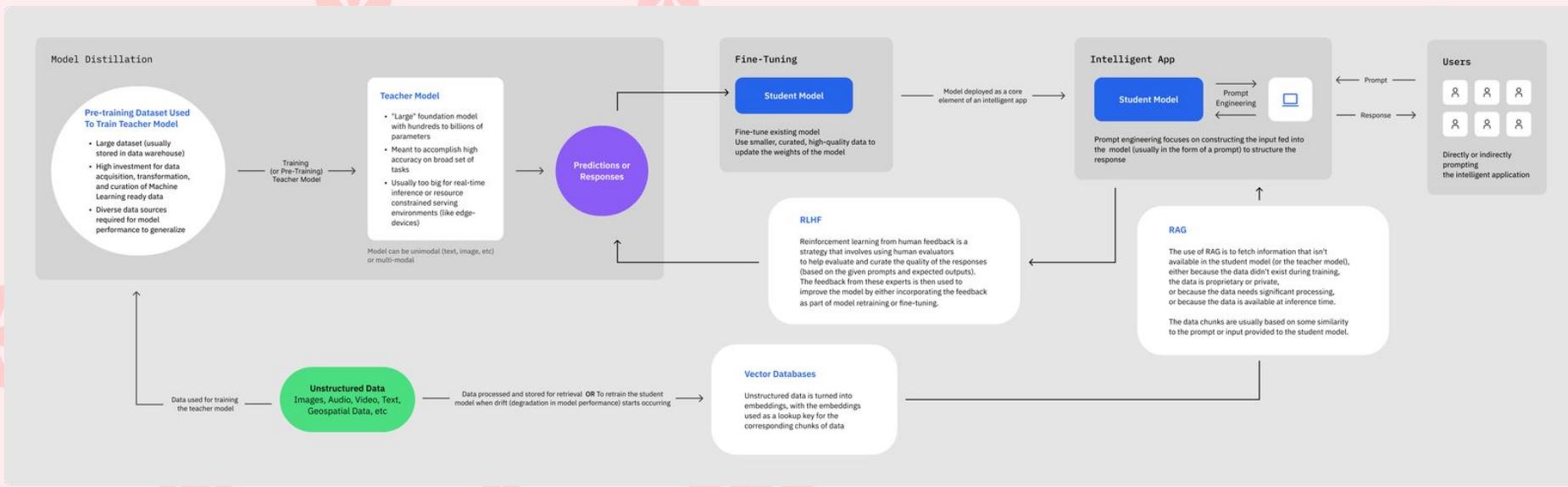
Production Gen-AI Applications are complicated & scaling requires considerations that are hard to address in the short-time we have



Discussion Question: Is ChatGPT an LLM? 🗣️



Production Gen-AI Apps: (Model Distillation +) Fine-Tuning + RAG + Evaluation _ RLHF

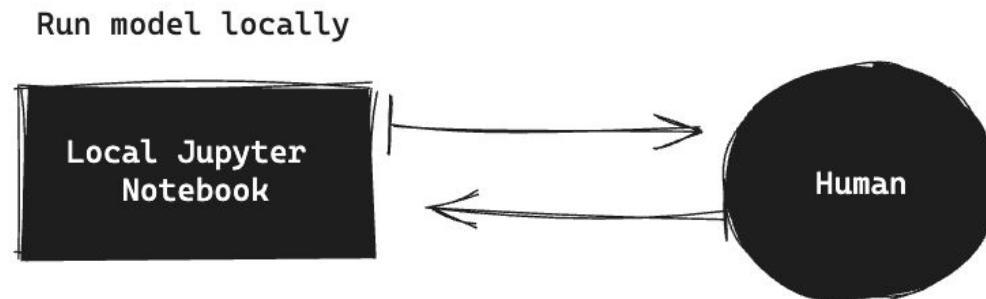


- Note: If you can't see the diagram, the original one is here:
<https://labelbox.com/blog/a-pragmatic-introduction-to-model-distillation-for-ai-developers/>

Easiest Way To Interact & Develop With An LLM

Most Simple Development

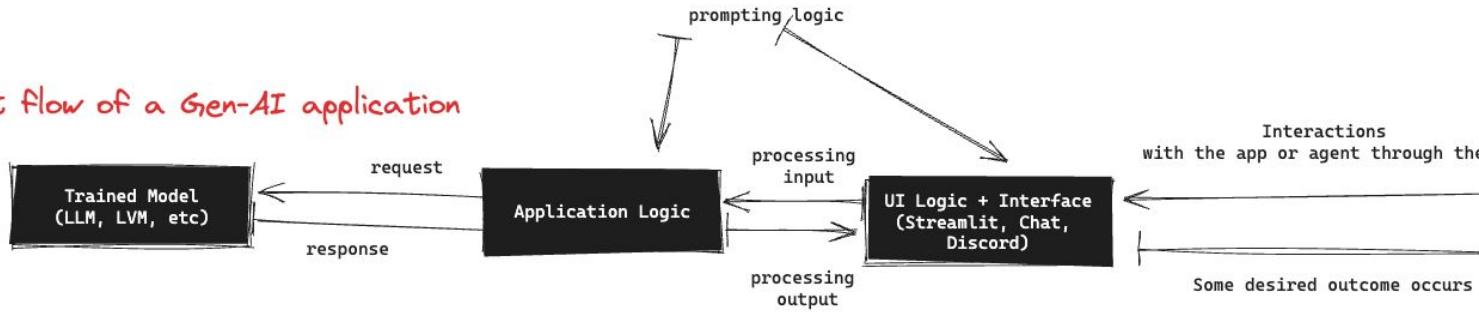
Most simple development using an OSS model



- You might use a tool like Ollama to run an LLM locally but for the most part you can be really slim on libraries like Langchain, etc
- UX will be pretty low-level & may need to specify “chat-like” prompts

Simplified Flow of Prompting

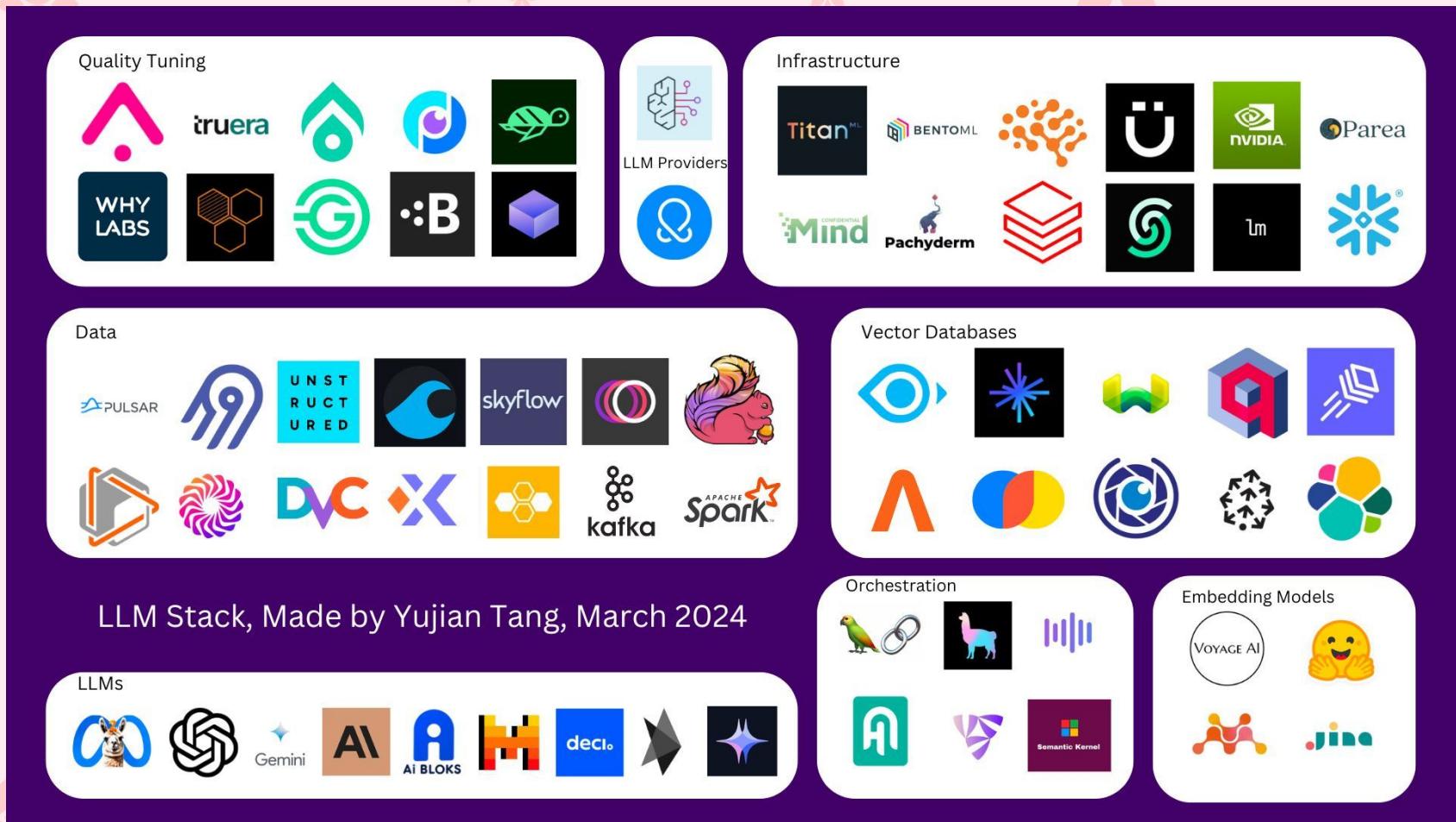
Simplest flow of a Gen-AI application



- OpenAI's Guide:

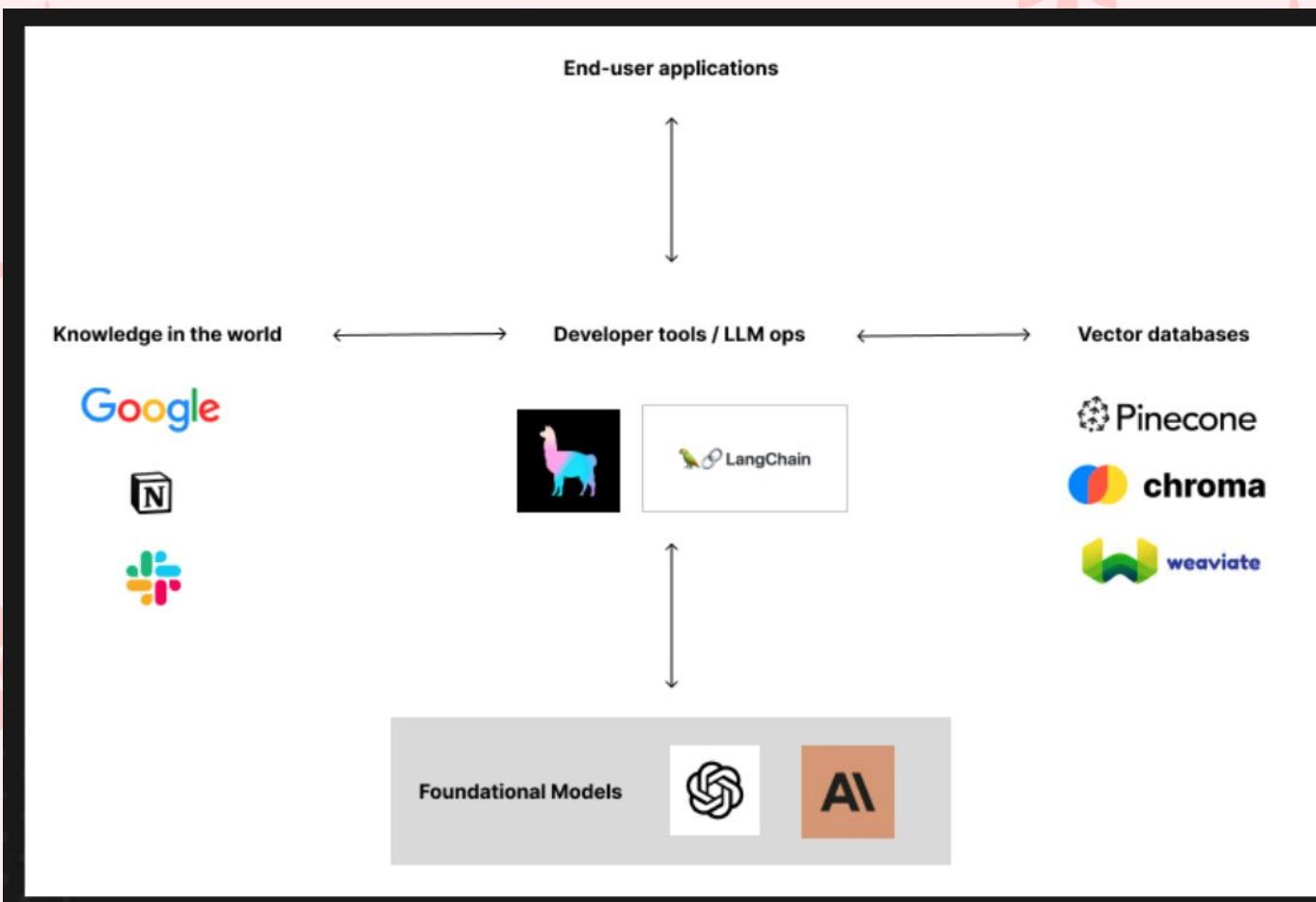
<https://platform.openai.com/docs/guides/prompt-engineering>

LLM Stack



- Medium:
<https://medium.com/plain-simple-software/the-lm-app-stack-2024-eac28b9dc1e7>

LLM Stack

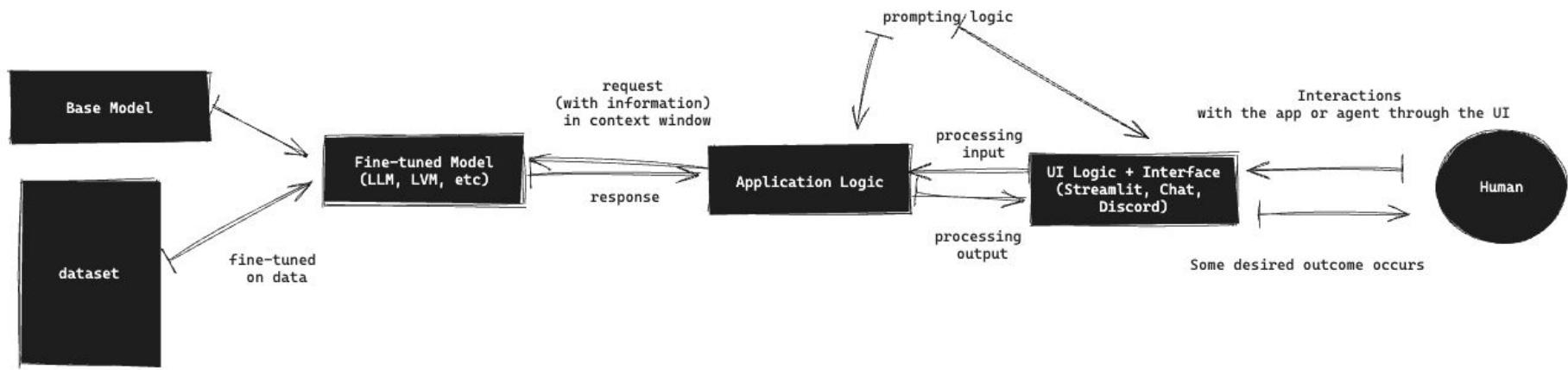


- AI Makerspace:

<https://www.canva.com/design/DAFsTCKwPGw/sxIXSvaU3C6HQCNRftp9g/view>

Simplified Flow of Fine-Tuned App

Flow of a Gen-AI application
with fine-tuned model



Why Open-Source? And who OS?



- Original link: [Choosing Your Path in Generative AI: Open-Source or Proprietary?](#)

● Original link:
[Choosing Your Path in Generative AI: Open-Source or Proprietary?](#)

Criteria	Commercial (Proprietary) Models	Open Source Models	Best Fit
Pricing	Usage-based, often with subscription models.	Generally free, potential costs for extra services.	Commercial: Ideal for businesses needing predictable costs. OSS: Fits budget-flexible or cost-sensitive projects.
Flexibility & Customization	Limited customization, standardized solutions.	High customization potential.	Commercial: Suited for standard solution requirements. OSS: Best for highly customized project needs.
Integration Capabilities	Better integration with existing ecosystems.	Requires more effort for integration.	Commercial: Fits businesses with existing infrastructure needing seamless integration. OSS: Suitable for environments where custom integration is practical.
Scalability	Designed for scalability and easier expansion.	Highly scalable but requires expertise.	Commercial: Ideal for rapidly growing businesses needing scalable solutions. OSS: Fits companies with the capability to manage scalability in-house.
Security & Compliance	Robust security and compliance with regulations.	Varies, may need additional work.	Commercial: Best for industries with stringent security and compliance needs. OSS: Suitable for scenarios, allowing for custom security adaptations.
Support & Maintenance	Structured support and maintenance.	Community-based, varies in quality.	Commercial: Suited for organizations needing consistent, dedicated support. OSS: Ideal for entities capable of self-maintenance and leveraging community support.
Community & Ecosystem	Certified partners and service providers.	Broad, collaborative community.	Commercial: Suitable for businesses seeking established service networks. OSS: Best for leveraging collaborative community input and development.
Long-term Viability	Predictable due to commercial backing.	Depends on community support.	Commercial: Ideal for organizations preferring stability and long-term support. OSS: Fits projects adaptable to community-driven changes and developments.

Considerations In Choosing Models

Accuracy
vs
Speed

Customizability
vs
Ease

Product
Requirements

Memory
vs
Computation

Generalized
vs
Specialized

Licensing

The De Facto Site For Models (Esp. OS)

The screenshot shows the Hugging Face website interface. At the top, there's a search bar with placeholder text "Search models, datasets, users...". Below the search bar, a navigation bar includes links for "Models", "Datasets", "Spaces", "Posts", "Docs", "Solutions", and "Pricing". On the far right of the navigation bar is a user profile icon.

The main content area is divided into two main sections: "Tasks" on the left and "Models" on the right.

Tasks: This section lists various model categories with their respective icons and names. Categories include:

- Multimodal: Image-Text-to-Text, Visual Question Answering, Document Question Answering
- Computer Vision: Depth Estimation, Image Classification, Object Detection, Image Segmentation, Text-to-Image, Image-to-Text, Image-to-Image, Image-to-Video, Unconditional Image Generation, Video Classification, Text-to-Video, Zero-Shot Image Classification, Mask Generation, Zero-Shot Object Detection, Text-to-3D, Image-to-3D, Image Feature Extraction
- Natural Language Processing: Text Classification, Token Classification, Table Question Answering, Question Answering, Zero-Shot Classification, Translation, Summarization, Feature Extraction, Text Generation, Text2Text Generation, Fill-Mask, Sentence Similarity
- Audio: Text-to-Speech, Text-to-Audio, Automatic Speech Recognition, Audio-to-Audio, Audio Classification, Voice Activity Detection
- Tabular: Tabular Classification, Tabular Regression
- Reinforcement Learning: Reinforcement Learning, Robotics
- Other: Graph Machine Learning

Models: This section displays a grid of 597,702 model cards. Each card contains the model name, a small profile picture, a brief description, and some statistics like the number of updates and the size of the model. Some examples of models listed include:

- CohereForAI/c4ai-command-r-plus
- mistral-community/Mistral-8x22B-v0.1
- mistralai/Mistral-7B-Instruct-v0.2
- CohereForAI/c4ai-command-r-plus-4bit
- ai21labs/Jamba-v0.1
- stabilityai/cosxl
- google/codgemma-7b-it
- Vezora/Mistral-22B-v0.1
- google/codegemma-7b
- cognitivecomputations/dolphin-2.8-mistral-7b-v02
- HuggingFaceH4/zephyr-orpo-141b-A35b-v0.1
- meta-llama/Llama-2-7b-chat
- ByteDance/SDXL-Lightning
- google/recurrentgemma-2b
- mistralai/Mistral-8x7B-Instruct-v0.1
- NexaAIDev/Octopus-v2
- jetmoe/jetmoe-8b
- google/gemma-1.1-7b-it
- v2ray/Mistral-8x22B-v0.1
- CohereForAI/c4ai-command-r-v01
- databricks/dbrx-instruct
- stabilityai/stablelm-2-12b
- dranger003/c4ai-command-r-plus-iMat.GGUF
- parler-tts/parler_tts_mini_v0.1
- qnguyen3/nanoLLaVA
- xai-org/grok-1
- Qwen/Qwen1.5-32B
- Qwen/Qwen1.5-32B-Chat
- google/recurrentgemma-2b-it
- openai/whisper-large-v3

At the bottom of the page, there are navigation links for "Previous" and "Next" with page numbers 1, 2, 3, and 4024.

The De Facto Site For Models (Esp. OS)

Multimodal

- Image-Text-to-Text
- Visual Question Answering
- Document Question Answering

Computer Vision

- Depth Estimation
- Image Classification
- Object Detection
- Image Segmentation
- Text-to-Image
- Image-to-Text
- Image-to-Image
- Image-to-Video
- Unconditional Image Generation
- Video Classification
- Text-to-Video
- Zero-Shot Image Classification
- Mask Generation
- Zero-Shot Object Detection
- Text-to-3D
- Image-to-3D
- Image Feature Extraction

Natural Language Processing

- Text Classification
- Token Classification
- Table Question Answering
- Question Answering
- Zero-Shot Classification
- Translation
- Summarization
- Feature Extraction
- Text Generation
- Text2Text Generation
- Fill-Mask
- Sentence Similarity

Audio

- Text-to-Speech
- Text-to-Audio
- Automatic Speech Recognition
- Audio-to-Audio
- Audio Classification
- Voice Activity Detection

Tabular

- Tabular Classification
- Tabular Regression

Reinforcement Learning

- Reinforcement Learning
- Robotics

Other

- Graph Machine Learning

Tasks Libraries Datasets Languages Licenses Other

Filter Datasets by name

- imagenet-1k
- mozilla-foundation/common_voice_7_0
- xtreme
- wikipedia
- mozilla-foundation/common_voice_11_0
- common_voice
- conll2003
- tweet_eval
- Open-Orca/OpenOrca
- marsyas/gtzan
- samsun
- bookcorpus
- fka/awesome-chatgpt-prompts
- LDJnr/Capybara
- clinc_oos
- OpenAssistant/oasst1
- intel/orca_dpo_pairs
- HuggingFaceH4/ultrafeedback_binarized
- c4
- kde4
- jondurbin/airobotos-2.2.1
- cnn_dailymail
- HuggingFaceH4/ultrachat_200k
- Open-Orca/SlimOrca
- garage-bAInd/Open-Platypus
- facebook/voxpopuli
- bigcode/starcoderdata
- mozilla-foundation/common_voice_13_0
- cerebras/SlimPajama-627B
- super_glue
- PolyAI/minds14
- ag_news
- google/fleurs
- teknum/openhermes
- databricks/databricks-dolly-15k
- billsum
- teknum/OpenHermes-2.5
- TIGER-Lab/MathInstruct
- librispeech_asr
- oscar
- huggan smithsonian_butterflies_subset
- universal_dependencies
- tiuae/falcon-refinedweb
- mc4
- migtissera/Synthia-v1.3
- Anthropic/hh-rlhf
- allenai/ultrafeedback_binarized_cleaned
- mozilla-foundation/common_voice_8_0
- togethercomputer/RedPajama-Data-1T
- meta-math/MetaMathQA
- pqiqa
- tatsu-lab/alpaca
- lmsys/lmsys-chat-1m
- jondurbin/truthy-dpo-v0.1

Tasks Libraries Datasets Languages Licenses Other

Filter Licenses by name

- apache-2.0
- mit
- openrail
- creativecommons-openrail-m
- other
- cc-by-nc-4.0
- llama2
- cc-by-4.0
- afl-3.0
- openrail++
- cc-by-nc-sa-4.0
- cc-by-sa-4.0
- gpl-3.0
- cc
- artistic-2.0
- bigscience-openrail-m
- bsd-3-clause
- bigscience-bloom-rail-1.0
- bigcode-openrail-m
- wtfpl
- cc-by-nc-nd-4.0
- cc0-1.0
- agpl-3.0
- cc-by-sa-3.0
- unlicense
- gpl
- bsd
- gemma
- cc-by-nc-2.0
- cc-by-3.0
- cc-by-2.0
- gpl-2.0
- bsd-2-clause
- lgpl-3.0
- bsl-1.0
- c-uda
- osl-3.0
- cc-by-nc-3.0
- cc-by-nd-4.0
- pddl
- ms-pl
- ecl-2.0
- zlib
- gfdl
- bsd-3-clause-clear
- cc-by-nc-sa-3.0
- lgpl
- mpl-2.0
- odbl
- deepfloyd-if-license
- cc-by-nc-sa-2.0
- lgpl-lr
- cc-by-nc-nd-3.0
- eupl-1.1
- cc-by-2.5
- odc-by
- epl-2.0
- isc
- cdla-permissive-2.0
- cdla-sharing-1.0
- ncsa
- etalab-2.0
- lgpl-2.1
- postgresql
- lppl-1.3c
- epl-1.0
- ofl-1.1
- cdla-permissive-1.0

Also has leaderboards: Open LLM Leaderboard

The screenshot shows the 'Open LLM Leaderboard' page on Hugging Face. At the top, there's a search bar and several filter buttons for 'Model types' (pretrained, continuously pretrained, fine-tuned on domain-specific datasets, chat models, base merges and moerges), 'Precision' (float16, bfloat16, 8bit, 4bit, GPTQ), and 'Model sizes' (in billions of parameters). The main area displays a table of models with columns for Average, ARC, HellaSwag, MMLU, TruthfulQA, Winogrande, and GSMBK. The table lists various models like 'davidkim205/Rhea-72b-v0.5', 'MTSAIR/MultiVerse_7B8', and 'SF-Foundation/Ein-72B-v0.11', along with their respective scores.

T	Model	Average	ARC	HellaSwag	MMLU	TruthfulQA	Winogrande	GSMBK
1	davidkim205/Rhea-72b-v0.5	81.22	79.78	91.15	77.95	74.5	87.85	76.12
2	MTSAIR/MultiVerse_7B8	81	78.67	89.77	78.22	75.18	87.53	76.65
3	MTSAIR/MultiVerse_7B8	80.98	78.58	89.74	78.27	75.09	87.37	76.8
4	SF-Foundation/Ein-72B-v0.11	80.81	76.79	89.62	77.2	79.62	84.06	78.77
5	SF-Foundation/Ein-72B-v0.13	80.79	76.19	89.44	77.07	77.82	84.93	79.3
6	SF-Foundation/Ein-72B-v0.12	80.72	76.19	89.46	77.17	77.78	84.45	79.23
7	abacusa1/Smaug-72B-v0.1	80.48	76.02	89.27	77.15	76.67	85.08	78.7
8	ibivibiv/Alpaca-Dragon-72B-v1	79.3	73.89	88.16	77.4	72.69	86.03	77.63
9	moezh/MoMo-72B-Lora-1.8.7-DPO	78.55	70.82	85.96	77.13	74.71	84.06	78.62
10	cloudyu/TomGrc_FusionNet_34Bx2_MoE_v0.1_DPO_f16	77.91	74.06	86.74	76.65	72.24	83.35	74.45
11	saltlux/luxia-21.4b-alignment-v1.0	77.74	77.47	91.88	68.1	79.17	87.45	62.4
12	cloudyu/TomGrc_FusionNet_34Bx2_MoE_v0.1_full_linear_DPO	77.52	74.06	86.67	76.69	71.32	83.43	72.93

- https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard

Also has leaderboards: LLM-Perf Leaderboard

The screenshot shows the LLM-Perf Leaderboard interface. At the top, there's a logo featuring a yellow emoji-like character with a crown and wings, surrounded by blue hexagons, with the word "Optimum" written below it. Below the logo is the title "(LLM-Perf Leaderboard)". A sub-section explains the purpose: "The LLM-Perf Leaderboard aims to benchmark the performance (latency, throughput, memory & energy) of Large Language Models (LLMs) with different hardwares, backends and optimizations using Optimum-Benchmark and Optimum flavors." It also states that anyone can request a model or configuration for automated benchmarking.

Below this, there are two main sections: "A100-80GB-275W" and "RTX4090-24GB-450W". Each section has a "About" button. Under each section, there's a "Control Panel" with various filters:

- Model:** A search bar with placeholder "Search for a model name".
- Open LLM Score (%):** A slider from 0 to 100.
- Peak Memory (MB):** A slider from 0 to 81920, currently at 81920.
- Backends:** A dropdown menu with "pytorch" selected.
- Load DTypes:** Checkboxes for "float32", "float16", and "bf16".
- Optimizations:** Checkboxes for "None", "BetterTransformer", and "FlashAttentionV2".
- Quantizations:** Checkboxes for "None", "BnB.4bit", "BnB.8bit", "GPTQ.4bit", "GPTQ.4bit+ExllamaV1", "GPTQ.4bit+ExllamaV2", "AWQ.4bit+GEMM", and "AWQ.4bit+GEMV".

At the bottom, there's a "Filter" button and a "Leaderboard" section with columns for Model, Arch, Params (B), Open LLM Score (%), Backend, Dtype, Opt., and Model ID. One row is visible: "cloudyai/Mixtral_110x2_MoE_19B" with Mixtral architecture, 19.19 params, 74.41 score, pytorch backend, bf16 dtype, and Non-optimal optimization.



- <https://huggingface.co/spaces/optimum/llm-perf-leaderboard>

Like TinyLlama

TinyLlama / **TinyLlama-1.1B-intermediate-step-1431k-3T** like 132

Text Generation Transformers PyTorch Safetensors cerebras/SlimPajama-627B bigcode/starcoderdata English llama Inference Endpoints text-generation-inference

License: apache-2.0

Model card Files and versions Community 8 Edit model card

TinyLlama-1.1B

<https://github.com/jzhang38/TinyLlama>

The TinyLlama project aims to **pretrain a 1.1B Llama model on 3 trillion tokens**. With some proper optimization, we can achieve this within a span of "just" 90 days using 16 A100-40G GPUs. The training has started on 2023-09-01.

We adopted exactly the same architecture and tokenizer as Llama 2. This means TinyLlama can be plugged and played in many open-source projects built upon Llama. Besides, TinyLlama is compact with only 1.1B parameters. This compactness allows it to cater to a multitude of applications demanding a restricted computation and memory footprint.

This Collection

This collection contains all checkpoints after the 1T fix. Branch name indicates the step and number of tokens seen.

Eval

Model	Pretrain								
	Tokens	HellaSwag	Obqa	WinoGrande	ARC_c	ARC_e	boolq	piaq	avg
Pythia-1.0B	300B	47.16	31.40	53.43	27.05	48.99	60.83	69.21	48.30
TinyLlama-1.1B-intermediate-step-50K-104b	103B	43.50	29.80	53.28	24.32	44.91	59.66	67.30	46.11
TinyLlama-1.1B-intermediate-step-500K-104b	503B	49.56	31.40	55.80	26.54	48.32	56.91	69.42	48.28

Downloads last month 91,999

Safetensors Model size 1.1B params Tensor type F32

Inference API Examples

My name is Mariama, my favorite

Compute +Enter 0.3

This model can be loaded on Inference API (serverless).

Model not loaded yet

JSON Output Maximize

Datasets used to train TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-3T

bigcode/starcoderdata Updated May 16, 2023 ± 4.85k 316

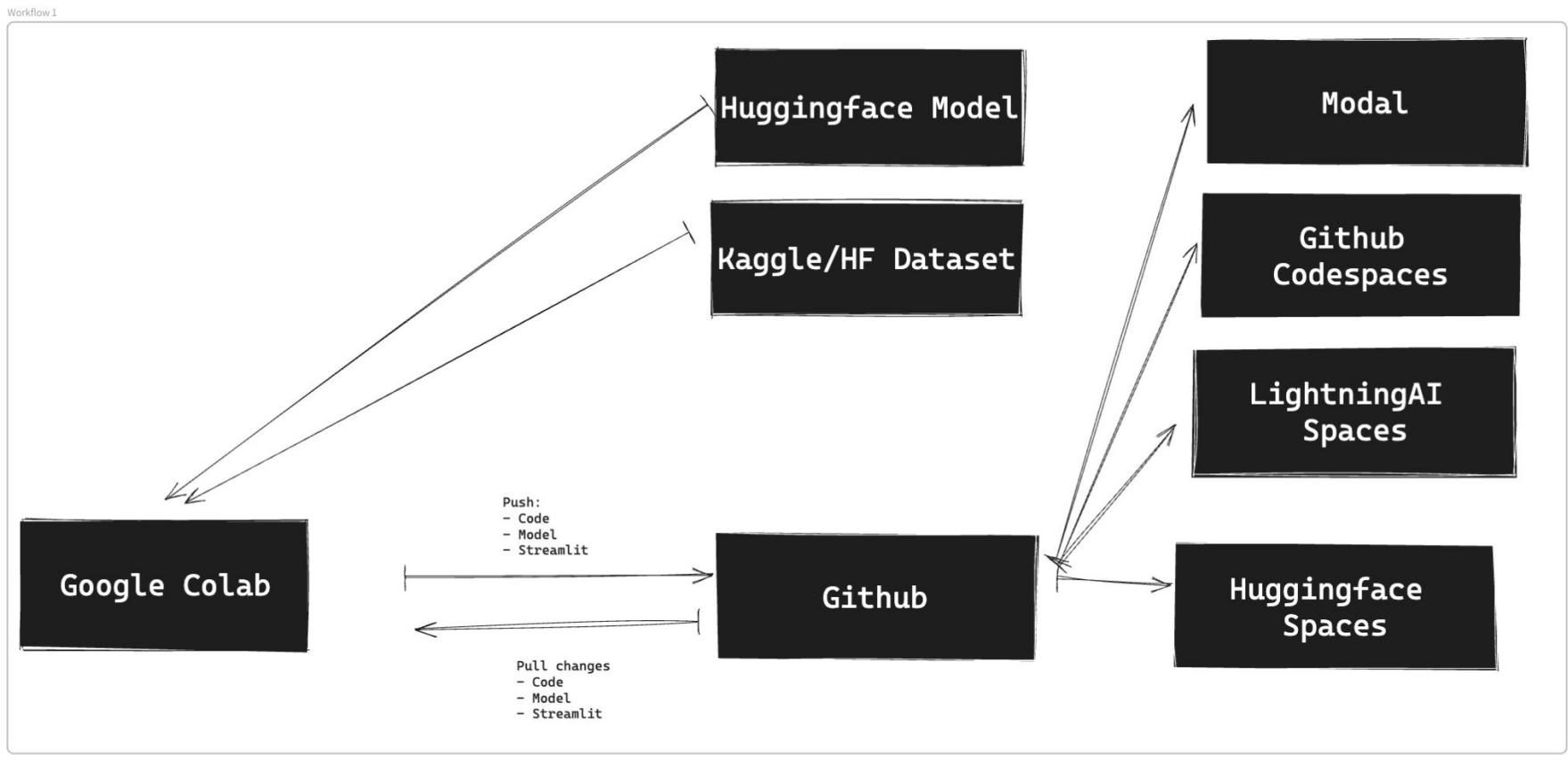
cerebras/SlimPajama-627B Updated Jul 7, 2023 ± 1.77k 336

Spaces using TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-3T 5

eduagarcia/open_pt_llm_leaderboard mhenrichsen/Axolotl_Launcher

View Details

Easiest Way To Interact & Develop With An LLM





Ex 1: Fine-Tuning a Taylor Swift TinyLlama Model

LO_GenAI_Workshops / oss-models / [Explainer]_HelloTaylorSwift_FineTuning.ipynb [🔗](#)

MMBazel Rename [Explainer]_HelloTaylorSwift_FineTuning.ipynb to oss-models/[E... [...](#)

Preview Code Blame 2717 lines (2717 loc) · 128 KB

[Open In Colab](#)

Fine-Tuning a TinyLlama (tinyllama_tayswifty) Model For Fun & Profit

TinyLlama is a 1.1B Llama model that is currently being trained on 3 trillion tokens, which recently started on September 1st. In this project, I fine-tune the latest version of TinyLlama to generate song lyrics in the style of Taylor Swift.

Source Materials

HelloTaylorSwift tutorial is based primarily on this tutorial:

- Original kaggle notebook
- Kaggle dataset

However these other resources are also helpful:

- Similar Model: <https://huggingface.co/huggingartists/taylor-swift>
- Similar Dataset: <https://huggingface.co/datasets/huggingartists/taylor-swift>

Tutorials on SFT & fine-tuning, TinyLlama, & HuggingFace

- Fine-Tune Your Own Tiny-Llama on Custom Dataset
- TinyLlama LLM: A Step-by-Step Guide to Implementing the 1.1B Model on Google Colab
- Instruct-Tune Llama to Create ChatGPT Like Chatbots | Custom Dataset, Huggingface, SFT
- <https://github.com/uygarkurt/SFT-TinyLlama/tree/main>

Tools

You'll need an access & accounts for:

- Google Colab - Ideally Pro (it's just faster to use a GPU like A100 or V100 High-RAM ~\$10)
- Huggingface - Also ideally Pro (there are some great benefits, including unlimited model and dataset upload ~ \$9)



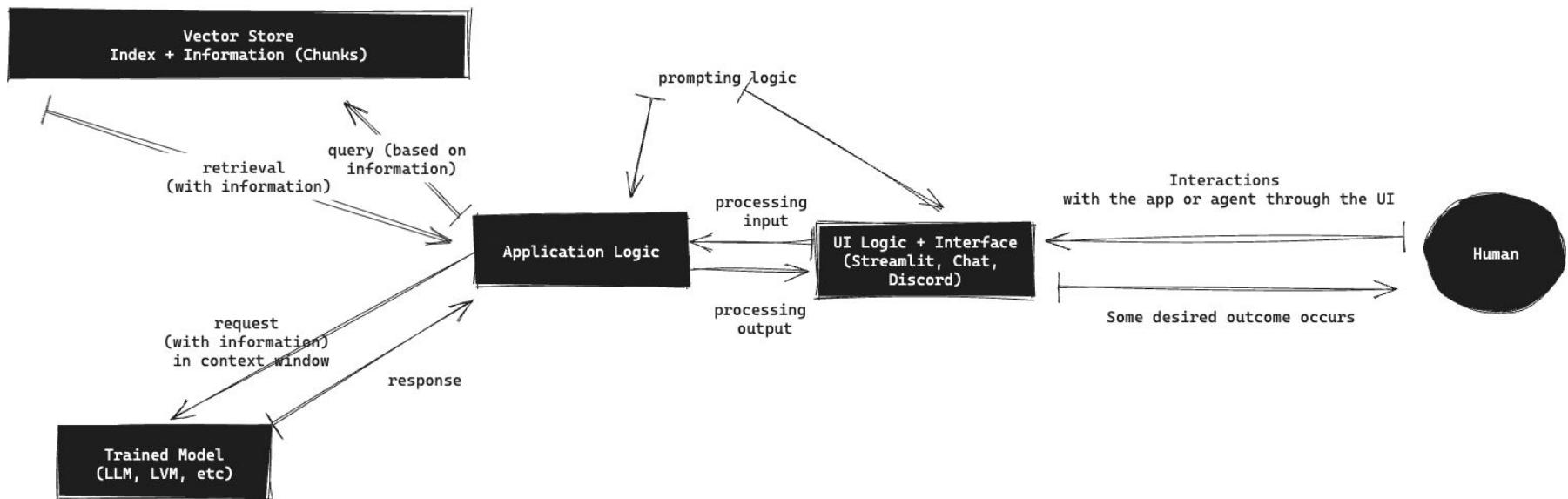
Ex 1: Fine-Tuning a Taylor Swift TinyLlama Model



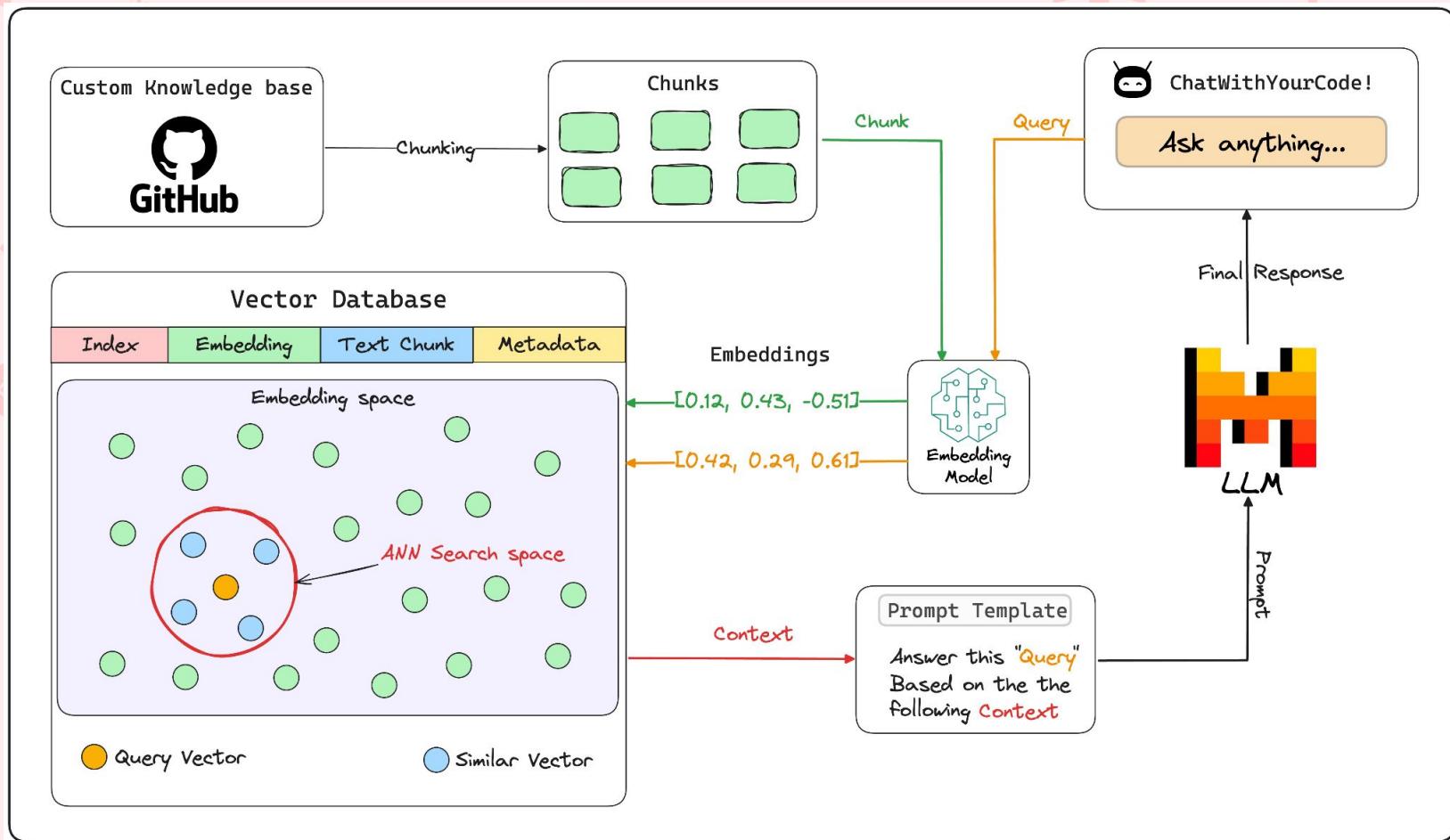
- [Explainer]_HelloTaylorSwift_FineTuning.ipynb:
<https://colab.research.google.com/drive/1W3HQDpU2kOJSq2FF8Qts1m110XopiQa4?usp=sharing>
- [Mini]_HelloTaylorSwift_FineTuning.ipynb:
<https://colab.research.google.com/drive/1hGKgfb0SmFeJDC12dHb4X7FqAu0vt6Oh?usp=sharing>

Simplified Flow of RAG

Flow of a Gen-AI application
with data (aka RAG)



Ex: Simple RAG



- Link: <https://lightning.ai/lightning-ai/studios/chat-with-your-code-using-rag?section=featured>

Ex: Simple RAG

Chat with your code using RAG!

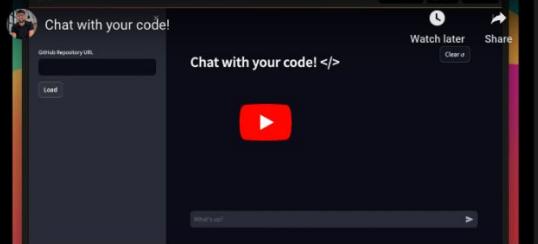
Akshay Pachaar · Published March 1, 2024

Overview Files

Chat with your code using RAG!

Transform your interaction with GitHub repositories through a natural language interface. In this studio we are building a "Chat with your code" RAG application that simplifies code queries, making coding more intuitive and productive.

And here's a quick demo of what we're building:



Chat with your code demo!

Open in Studio

Chat With your code! </>
Powered by LlamaIndex!

Machine: (1 x A10G) GPU
License: Apache-2.0
Get Studio badge

Text Tutorials Blogs

Chat with your code using RAG!

Try it yourself
Run main notebook
Chat with your code app!

Key architecture components

- Custom knowledge base
- Chunking
- Embeddings model
- Vector databases
- User chat interface
- Query engine
- Prompt template

Conclusion
Next steps

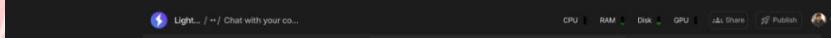
Try it yourself

At this point I would encourage you to take it for a spin by clicking on the [Open Template](#) & follow the steps below:

The studio launches in a new tab, you can always come back here to understand how things work!

Run main notebook

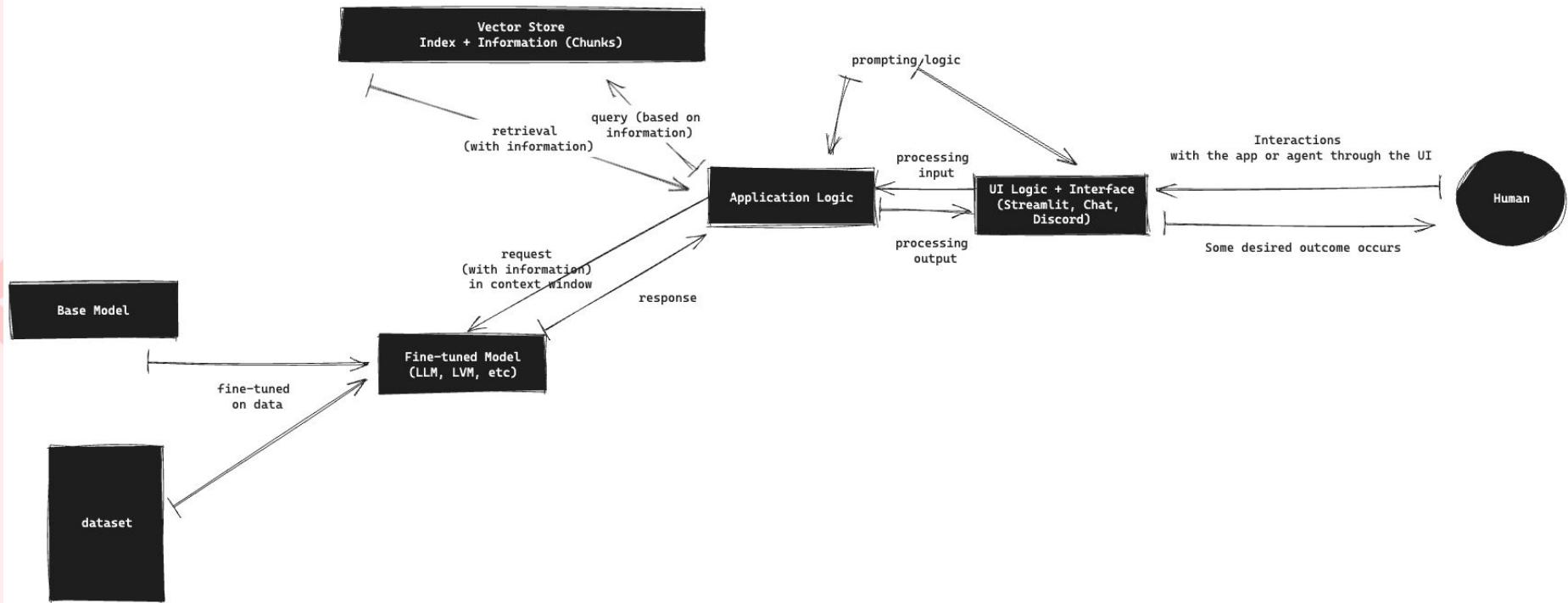
You can start by running the [main.ipynb](#) notebook, which contains the essential code to set up a query engine for interacting with the repository you provide.



- Link: <https://lightning.ai/lightning-ai/studios/chat-with-your-code-using-rag?section=featured>

Simplified Flow of Rag + Fine-Tuned Based Application

Flow of a Gen-AI application
with data (aka RAG) & fine-tuned model



Prompting vs Fine-Tuning vs RAG

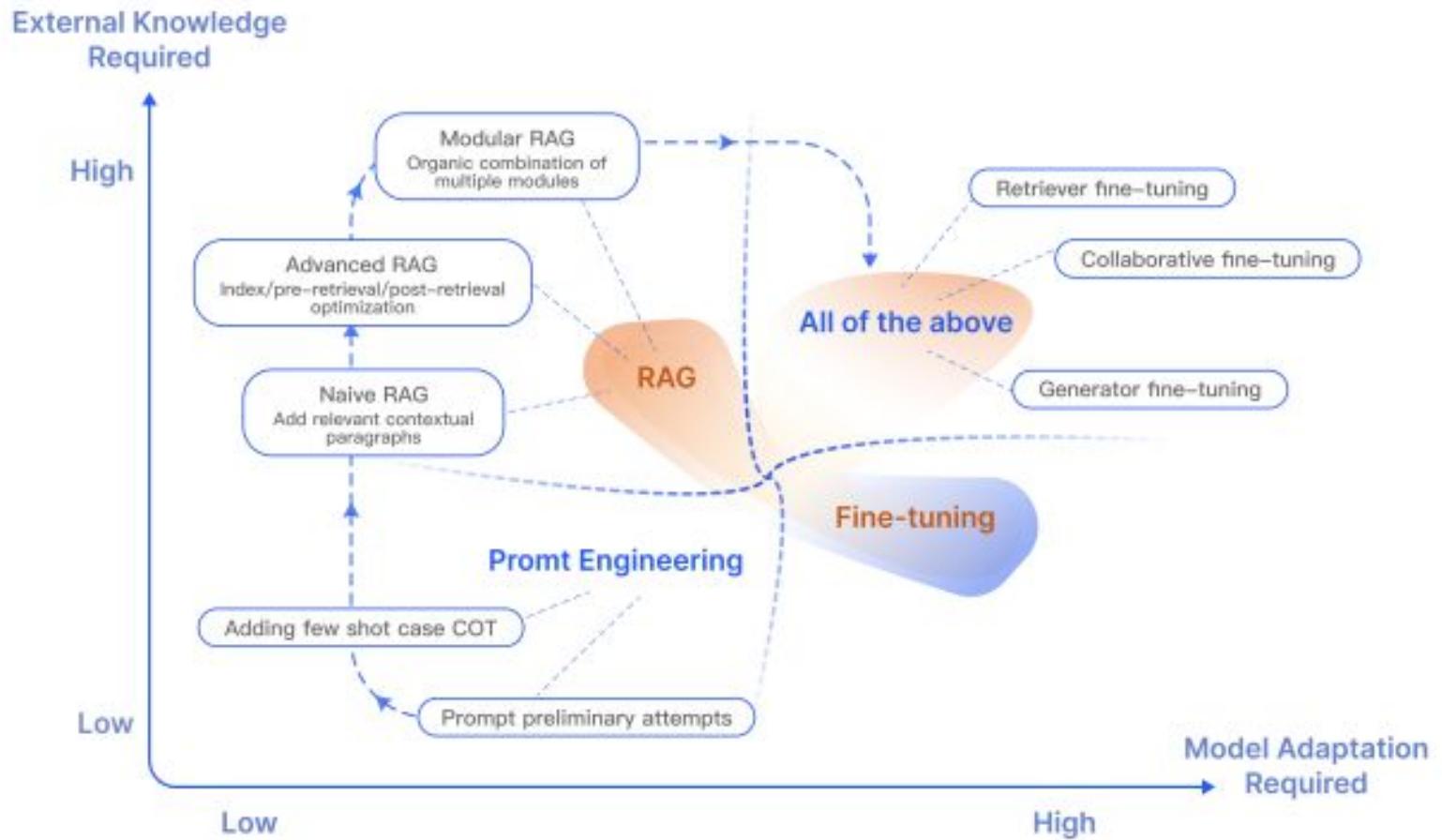


Figure 2: RAG compared with other model optimization methods

- Note: Comparison of RAG versus Fine-Tuning – Source: "[Retrieval-Augmented Generation for Large Language Models: A Survey](#)"



Ex 2: RAG With A Vector Database



Chat with your code: RAG with Weaviate and Llamaindex

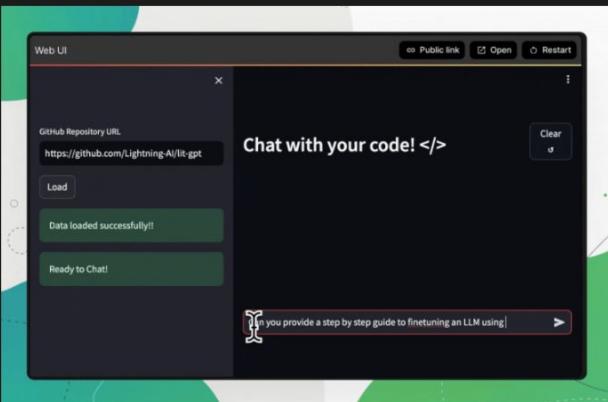
Leonie Published April 5, 2024

Overview Files

Open In Studio

Chat with your code: RAG with Weaviate and Llamaindex

This Studio template builds a simple application that lets you chat with code from a GitHub repository. Here's a quick demo of what we're building:

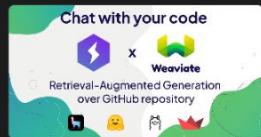


Demo of "Chat with your code" app.

Under the hood, you are building a Retrieval-Augmented Generation (RAG) application using the following components:

- **Embedding model:** BGE embedding model via [Hugging Face](#)
- **Vector database:** [Weaviate](#)
- **Large Language Model (LLM):** [Mistral](#) via [Qilama](#)
- **Orchestration framework:** [Llamaindex](#)
- **App framework:** [Streamlit](#)

Note: This Studio template is based on the Studio template "["Chat with your code using RAG!"](#)" by [Akshay Pachaaar](#).



Machine: (1 x A10G) GPU

License: Apache-2.0

Get Studio badge

[Text](#) [Tutorials](#)

Chat with your code: RAG with Weaviate and Llamaindex

Prerequisites

- Set up a Weaviate vector database
- Set up environment variables

Build a Retrieval-Augmented Generation Pipeline

1. Setup LLM and embedding model
2. Load files from GitHub repository
3. Prepare Weaviate vector database
4. Setup query engine
5. Setup a simple user interface

Summary

- Link:

<https://lightning.ai/weaviate/studios/chat-with-your-code-rag-with-weaviate-and-llamaindex>

The Ops Part of LLMOps

KEY PIECES, LLM OPS DEV

Is it just MLOps + LLMs?

- e.g., GitHub, Conda, UNIX Terminal, VS Code, Jupyter Notebook + OpenAI API

Kind of

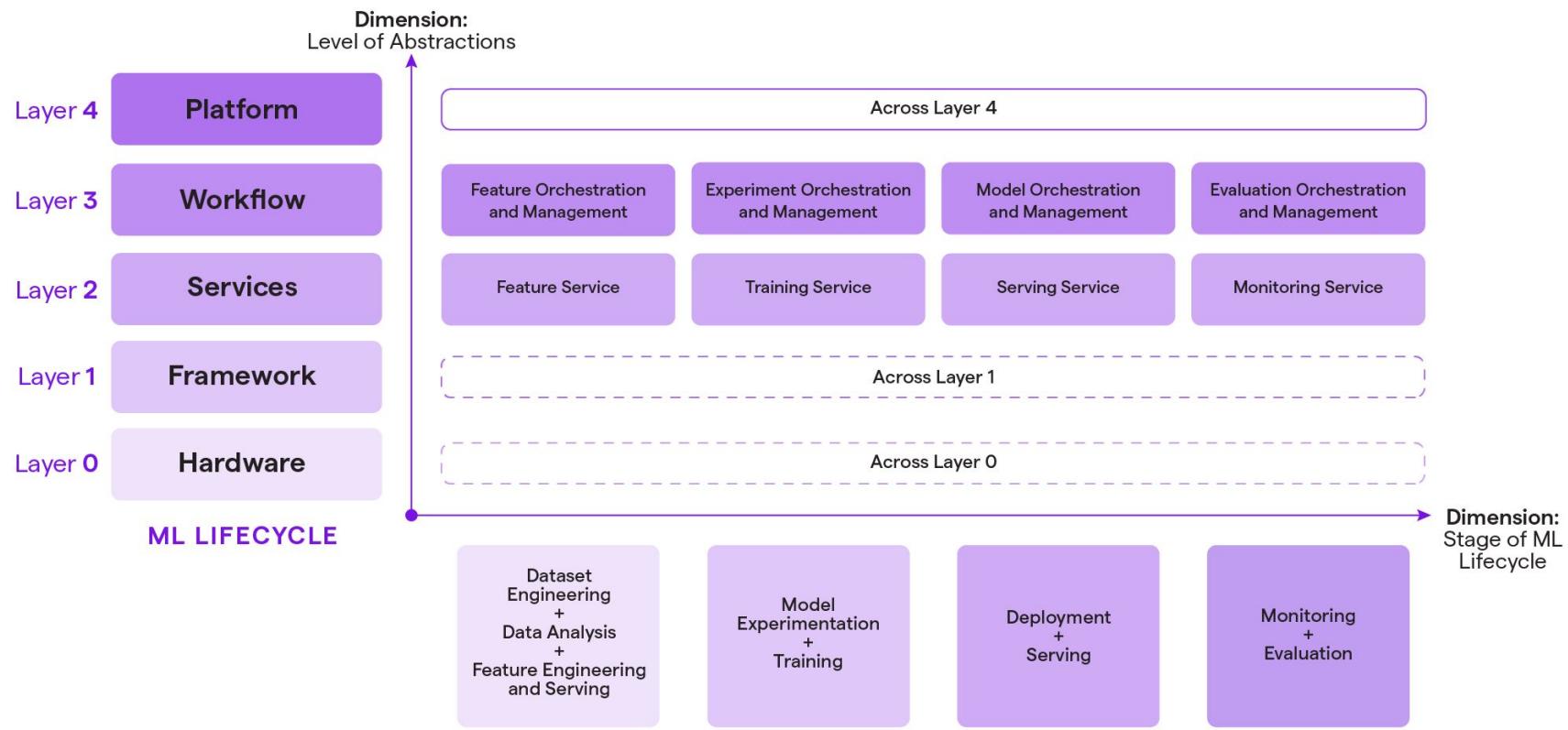
- But also, FastAPI + Docker => HF

And of course

- AWS, GCP, Azure, also Databricks...



MLOps Stack (Pre-GenAI)



- Link:
<https://www.featureform.com/post/measuring-your-ml-platforms-north-star-metrics>

MLOps Stack (Pre-GenAI)

The Four Wars of the AI Stack

The Quality Data Wars	
Part of the AI Stack Wars	
Content platforms defending from content raids	
Belligerents	
Journalists	AI researchers
Writers	AI startups
Artists	Synthetic Data research
Leaders and flagbearers	
The New York Times	OpenAI
StackOverflow	Stability AI
Reddit	Microsoft
Getty Images	CommonCrawl
Sarah Silverman	Axel Springer
Greg Rutkowski	X (Elon Musk)
Sarah Andersen	Eleuther AI
George R. R. Martin	LAION / DataComp

The War of the GPU Rich/Poors	
Part of the AI Stack Wars	
Faster/Cheaper Inference, Finetuning & Training	
Belligerents	
GPU Rich Clouds	GPU Poor AI Engineers
GPU Manufacturers	Edge/Local Compute
VC Funding	New Model Research
Leaders and flagbearers	
Nvidia	Modular / MLC
Google	TinyCorp
Microsoft	QLoRA (Tim Dettmers)
Amazon Bedrock	r/LocalLlama
Together.ai	Consistency Models
Fireworks.ai	Apple / MLX
Anyscale	RWKV / RetNet
Replicate	Mamba / StripedHyena

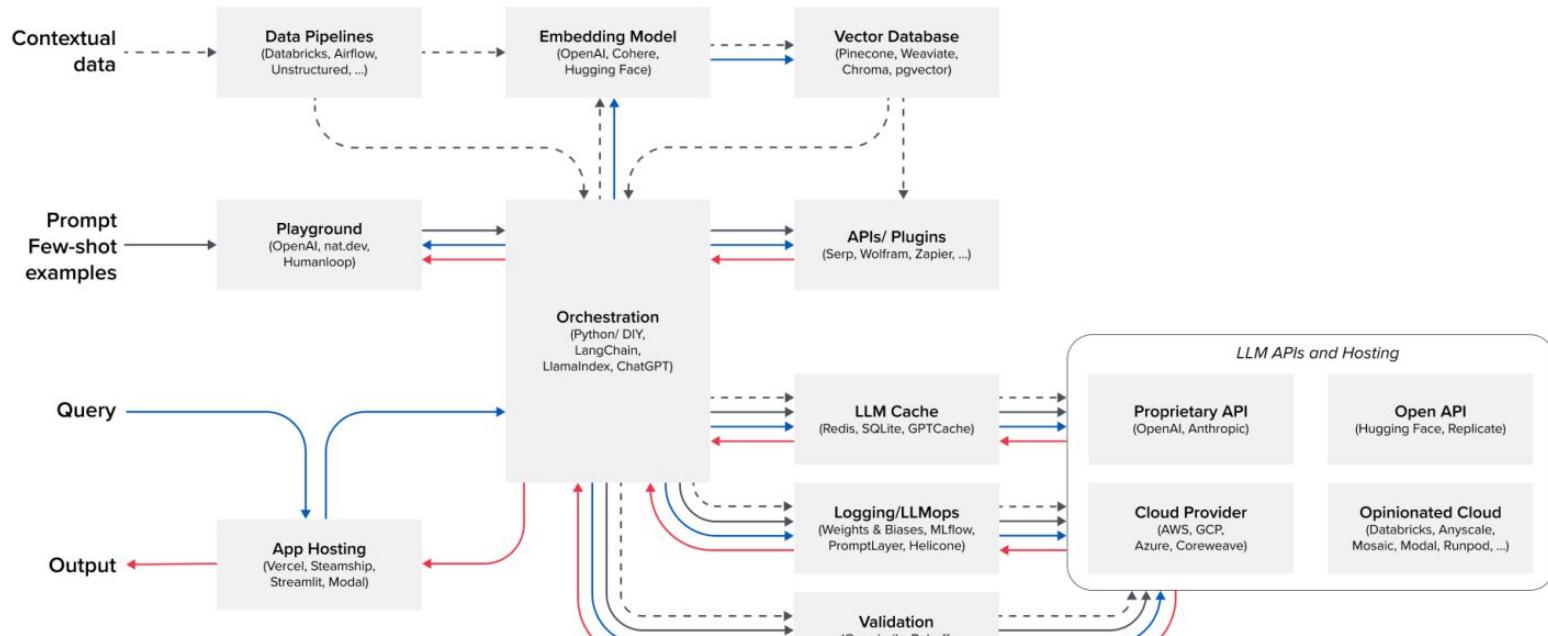
The Multimodality War	
Part of the AI Stack Wars	
Specialist models vs Everything models	
Belligerents	
Text to Image Startups	OpenAI
Text to Audio Startups	Google Deepmind
Other modalities	
Leaders and flagbearers	
Midjourney	OpenAI ChatGPT/API
Playground.ai	Google Gemini
Lexica.art	
Eleven Labs	
Suno.ai	
Pika Labs	
HeyGen	
Assembly.ai	

The RAG/Ops War	
Part of the AI Stack Wars	
Databases, Frameworks, vs Dev Tooling	
Belligerents	
General DBs	LLM Frameworks
Vector DBs	LLM Platforms
Search/Recsys	AIEF Standards
Leaders and flagbearers	
Postgres (pgvector)	LangChain
MongoDB Vector	LlamaIndex
Cassandra Vector	Guardrails
Elasticsearch	Microsoft Autogen
Redis Vector	Flowise
Pinecone	BuildShip
Chroma	Vellum
Qdrant	LangServe
	LangSmith
	Gantry

- Link: Latent.Space: <https://www.latent.space/p/dec-2023-audio>

Gen-AI Stack???

Emerging LLM App Stack



Gen-AI Stack???

Data pipelines	Embedding model	Vector database	Playground	Orchestration	APIs/plugins	LLM cache
Databricks	OpenAI	Pinecone	OpenAI	Langchain	Serp	Redis
Airflow	Cohere	Weaviate	nat.dev	Llamaindex	Wolfram	SQLite
Unstructured	Hugging Face	ChromaDB	Humanloop	ChatGPT	Zapier	GPTCache
		pgvector				

Logging / LLMops	Validation	App hosting	LLM APIs (proprietary)	LLM APIs (open)	Cloud providers	Opinionated clouds
Weights & Biases	Guardrails	Vercel	OpenAI	Hugging Face	AWS	Databricks
MLflow	Rebuff	Steamship	Anthropic	Replicate	GCP	Anyscale
PromptLayer	Microsoft Guidance	Streamlit			Azure	Mosaic
Helicone	LMQL	Modal			CoreWeave	Modal
						RunPod



Have Additional Q's For Me About MLOps? Find Me On



[in/mikikobazeley](https://www.linkedin.com/in/mikikobazeley)



[BazeleyMikiko](https://twitter.com/BazeleyMikiko)



[MikiBazeleyTheMLOps
Engineer](https://www.youtube.com/c/MikiBazeleyTheMLOpsEngineer)



[Click Long Link](#)