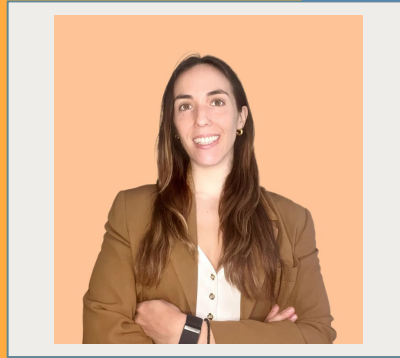# Life Beyond ChatGPT

From RAG to fine tuned models with enterprise data

# Presenter

Maria Zervou

Sr. Specialist Solutions Architect

# Agenda

- New Wave of Deep Learning
- Customisation Phases of GenAI
- Build a RAG application
- Live demo
- Fine Tuning Concepts
- 5 mins on Pretraining
- Do you really want to discuss cost?
- Summary – Call to Action

# Look out for....

# Look out for….



Expert

Newbie

Extra HOt

HOt

Medium

Mild

Extra Mild

PERi-Ometer

# Agenda

- **New Wave of Deep Learning**
- Customisation Phases of GenAI
- Build a RAG application
- Live demo
- Fine Tuning Concepts
- 5 mins on Pretraining
- Do you really want to discuss cost?
- Summary – Call to Action

# What are GenAi Models

**Large Language Models**



**Diffusion Models**

# What are GenAi Models

LARGE AI MODELS

To Train Generative AI models we need…..

GPUs

**How many do you think we need?**

To Train Generative AI models we need…..

> # **GPUs**

**We need a lot of GPUS to train
your own Generative AI models**

# Transformer Models

Fueling the next wave of Deep Learning



2017

Large Language Models
(Transformer)

Applications of Transformer Models

Transformer

Recommendation Agents

Image Classification

Conventional Models
Train on **Labeled** Datasets

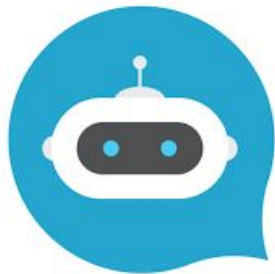Transformer Models
Train on **Unlabeled** Datasets

# How Gen Ai will disrupt you?
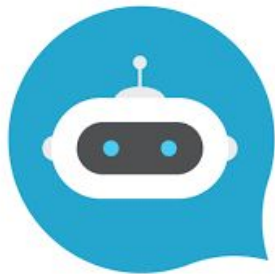
# How it will disrupt you?

Data + GenAi = Huge business value



**Create Conversational Interfaces
for everything**

# How it will disrupt you?
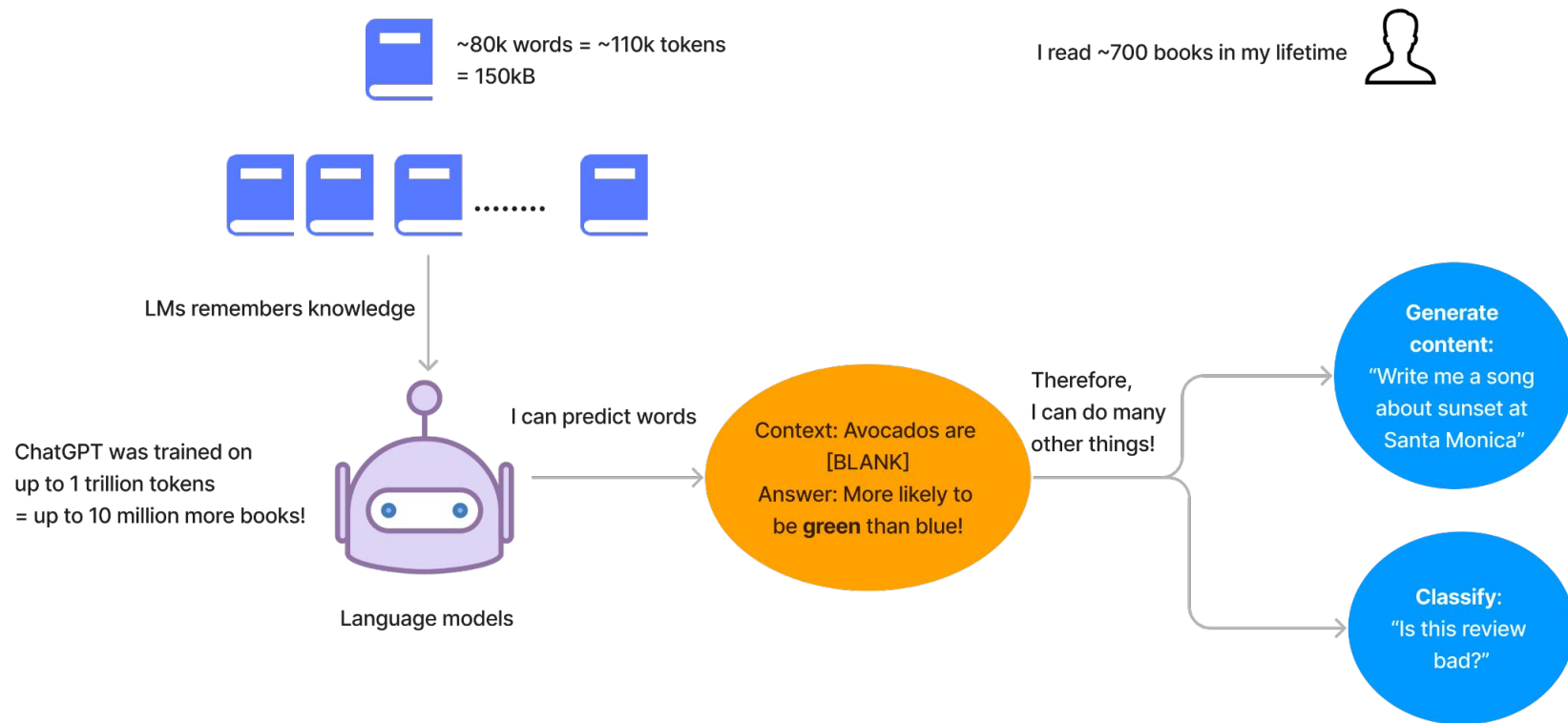
Data + GenAi = Huge business value



**Create Conversational Interfaces for everything**



**Human Level Comprehension but faster #Notjustachatbot**

# Human Level Comprehension

~80k words = ~110k tokens = 150kB

I read ~700 books in my lifetime

LMs remembers knowledge

ChatGPT was trained on up to 1 trillion tokens = up to 10 million more books!

Language models

I can predict words

Context: Avocados are [BLANK]
Answer: More likely to be **green** than blue!

Therefore, I can do many other things!

Generate content: "Write me a song about sunset at Santa Monica"

Classify: "Is this review bad?"

# How it will disrupt you?

Data + GenAi = Huge business value



**Create Conversational Interfaces for everything**



**Human Level Comprehension but faster #Notjustachatbot**



**Generate Human Quality Text and Images**

# So where do we start?

# Agenda

- New Wave of Deep Learning
- **Customisation Phases of GenAI**
- Build a quick RAG application with your own data
- Live demo
- Fine Tuning Concepts
- 5 mins on Pretraining
- Do you really want to discuss cost?
- Summary – Call to Action

# Customisation Phases of GenAi

**Foundational model with Prompts**

**Retrieval Augmented knowledge (RAG)**

**Fine-tune foundational model on your data**

**Fully retrain foundational models (your own "GPT")**

# Customisation Phases of GenAi

Mild

Extra Hot

| |
|---|
| **Foundational model with Prompts** |
| **Retrieval Augmented knowledge (RAG)** |
| **Fine-tune foundational model on your data** |
| **Fully retrain foundational models (your own "GPT")** |

# Kpis Time

| |
|---|
| **Foundational model with Prompts** |
| **Retrieval Augmented knowledge (RAG)** |
| **Fine-tune foundational model on your data** |
| **Fully retrain foundational models (your own "GPT")** |

**Training Cost**

**Data Size**

**Know-how**

**Customisation**

# Kpis Time

| |
|---|
| **Foundational model with Prompts** |
| **Retrieval Augmented knowledge (RAG)** |
| **Fine-tune foundational model on your data** |
| **Fully retrain foundational models (your own "GPT")** |

**Training Cost**

**Data Size**

Extra Hot

**Know-how**

**High**

**Customisation**

# Proprietary LLMs



**Proprietary LLMs**

**Pros:**
- Access to state-of-the-art models (e.g. GPT4).
- Easy to use: No need for advanced knowledge on LLM.

# Proprietary LLMs



**Proprietary LLMs**

**Pros:**
- Access to state-of-the-art models (e.g. GPT4)
- Easy to use: No need for advanced knowledge on LLM.

**Cons:**
- Model updates and service SLA depend fully on the providers.
- Quality of inference can vary significantly over time [ref].
- You **may** need to send your data to the providers (but not always).

# Open Source LLMs



**stability.ai**
**Stable Diffusion**

**Open Source LLMs**

**mosaic^ML**
**MPT**

🤗 **Hugging Face**

**databricks**
**Dolly**

Open source models customers can get from an open repository like Hugging Face.

**Pros:**
- Access to many general & specialized models.
- Data and models stay within your environment.
- Flexibility in tuning latency and throughput by changing the inference cluster configuration
- Transparency in the source code, model weights, and training dataset
- Fine-tuning is possible

# Open Source LLMs

Open source models customers can get from an open repository like Hugging Face.

**Pros:**
- Access to many general & specialized models.
- Data and models stay within the customers' environment.
- Flexibility in tuning latency and throughput by changing the inference cluster configuration
- Transparency in the source code, model weights, and training dataset
- Fine-tuning is possible

**Cons:**
- Requires expertise for tuning models and selecting the right infrastructure.

# Foundational Models

| |
|---|
| **Foundational model with Prompts** |
| **Retrieval Augmented knowledge (RAG)** |
| **Fine-tune foundational model on your data** |
| **Fully retrain foundational models (your own "GPT")** |

# Prompt Engineering

# Kpis Time – Foundational Models

**Foundational model with Prompts**

**Retrieval Augmented knowledge (RAG)**

**Fine-tune foundational model on your data**

**Fully retrain foundational models (your own "GPT")**

**Training Cost**

**Data Size**

**Extra Hot**

**Know-how**

**High**

**Customisation**

# Kpis Time – Foundational Models



| | Training Cost | Data Size | Know-how | Customisation |
|---|---|---|---|---|
| **Foundational model with Prompts** | 💰 ↓ | 🗄 ↓ | Mild ↓ | Low ↓ |
| **Retrieval Augmented knowledge (RAG)** | | | | |
| **Fine-tune foundational model on your data** | | | | |
| **Fully retrain foundational models (your own "GPT")** | 💰💰💰 | 🗄🗄 | Extra Hot | High |

DON'T FORGET!

# Customisation Phases of GenAi

Saas with Prompts

| | |
|---|---|
| **Foundational model with Prompts** | - **Ready to be plugged** to your applications<br>- Can modify your output with your own **Prompts**<br>- Can **generate "any answer"**<br>- - - - - - - - - - - - - - - - - - - - - - - - - -<br>- You **don't control data inside** the model knowledge base<br>- You **cannot control the model and version**<br>- You **cannot control ownership** |

DON'T FORGET!

# Customisation Phases of GenAi – RAG

Foundational model as SaaS with Prompts

Retrieval Augmented knowledge (RAG)

**Provide your data as** you are calling the model

Fine-tune foundational model on your data

Fully retrain foundational models (your own "GPT")

# RAG Architecture

# Agenda

- New Wave of Deep Learning
- Customisation Phases of GenAI
- **Build a quick RAG application with your own data**
- Fine Tuning Concepts
- Live demo
- 5 mins on Pretraining
- Do you really want to discuss cost?
- Summary – Call to Action

# RAG Architecture for Q/A Bot

On Insurance Policy Data



Insurance Policy pdf

txt

Vector Store

Embedded Question

Question + Context

Prompt

Response

LLM

question

answer

# Document Indexing



Policy Schedule / Validation Certificate

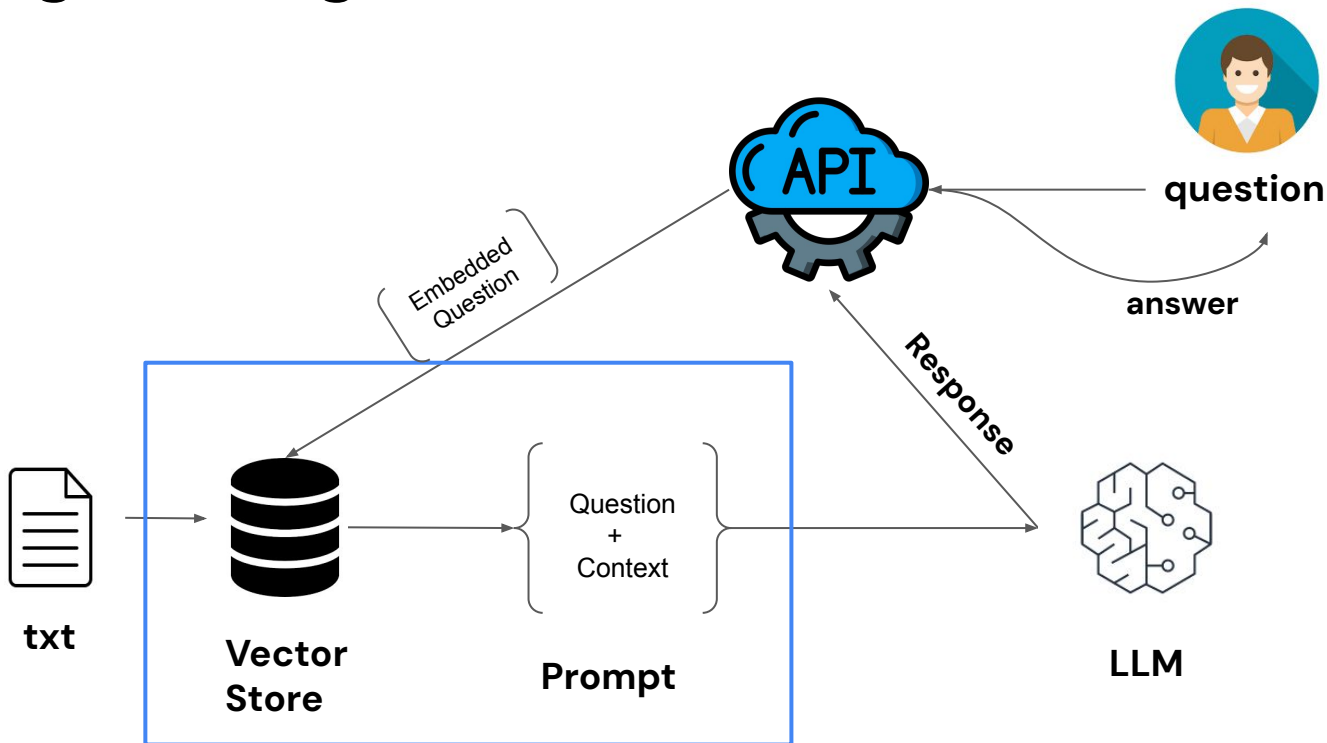Darth Vader
54 Death Star, Nebula
London

(Please attach to policy document)
SINGLE CAR HIRE EXCESS INSURANCE POLICY
Certification          No. BICEWCARTK/960
Underwriter:          Newline Insurance Company Limited
Issue Date:          16/07/2021
Start Date:          07/08/2021 Time: 11:00 AM
DropOff Date:       14/08/2021 Time: 5:30 PM
Total Days           8 Days
Area                   Europe
Lead Driver First Name    Darth
Lead Driver Surname      Vader
Address              54 Death Star, Nebula
                        London
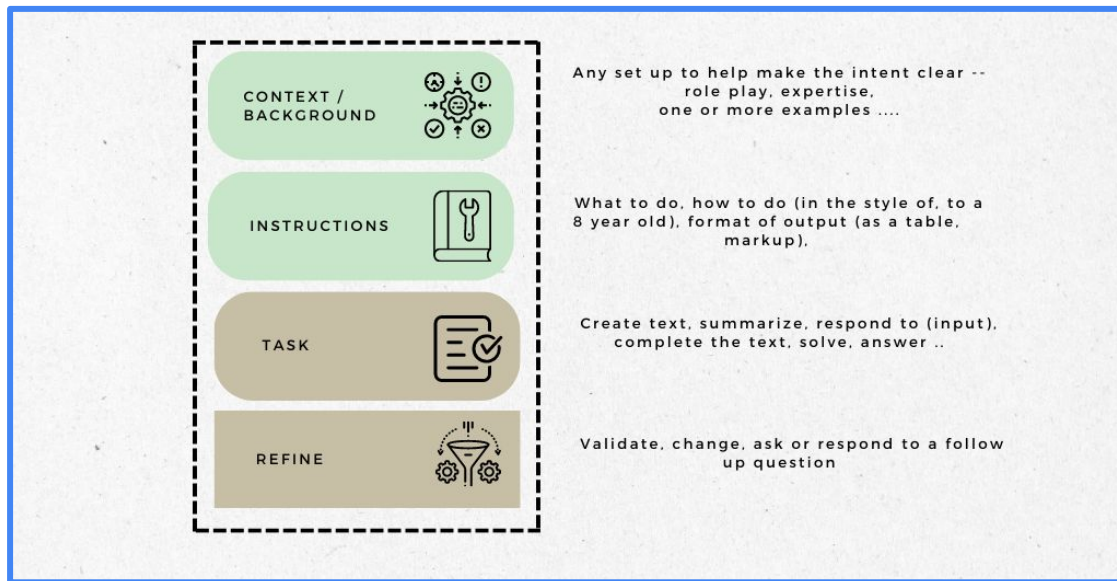Email                darthvader@yahoo.com
Mobile Number        791707777

| First Name | Surname | Age |
|------------|---------|-----|
| Darth | Vader | 245 |

Insurance Price: £17.86

**Insurance Policy pdf**

**txt**

**Vector Store**

Question + Context

**Prompt**

**LLM**

Embedding

Response

**question**

**answer**

# Document Indexing

# Embeddings

## Translate text to fixed size vectors



- Sentence vs Word Embeddings
- Chunk size
- Chunk strategy (Fixed , Recursive ...)
- Chunk Overlap
- Prompt Limits

Search query

Embedding Generator

Documents
(products, users, text)

Embedding Generator
(API, Transformer Model etc)

Dense Vectors

Vector Store

ENN or ANN search

Search Results

# Vector Index

# MVP Steps – Vector Store

# Prompt Engineering



**Insurance Policy pdf**

**txt**

**Vector Store**

Question + Context

**Prompt**

Embedded Question

Response

**question**

**answer**

**LLM**

# Prompt Engineering – Anatomy of a Prompt

# Prompt Engineering– Anatomy of a Prompt

You are a helpful assistant and you are good at helping to answer a question based on the context provided, the context is a document.

If the context does not provide enough relevant information to determine the answer, just say I don't know. If the context is irrelevant to the question, just say I don't know. If you did not find a good answer from the context, just say I don't know. If the query doesn't form a complete question, just say I don't know.

If there is a good answer from the context, try to summarize the context to answer the question.

The **"act as"** hack

# Prompt Engineering – Anatomy of a Prompt

You are a helpful assistant and you are good at helping to answer a question based on the context provided, the context is a document.

If the context does not provide enough relevant information to determine the answer, just say I don't know. If the context is irrelevant to the question, just say I don't know. If you did not find a good answer from the context, just say I don't know. If the query doesn't form a complete question, just say I don't know.

If there is a good answer from the context, try to summarize the context to answer the question.

Instructions

# Prompt Engineering – Anatomy of a Prompt

You are a helpful assistant built by Databricks, you are good at helping to answer a question based on the context provided, the context is a document.

If the context does not provide enough relevant information to determine the answer, just say I don't know. If the context is irrelevant to the question, just say I don't know. If you did not find a good answer from the context, just say I don't know. If the query doesn't form a complete question, just say I don't know.

If there is a good answer from the context, try to summarize the context to answer the question.

Task:
**Summarisation**

# Prompt engineering patterns

"Add Context to the Query"

Few shot prompting

User supplied examples {

Query
Example 1
Example 2
...

**Prompts can be:**

Natural language sentences or questions. Code snippets or commands. Combinations of the above.

Emojis....basically any text!

Prompts can include outputs from other LLM queries.
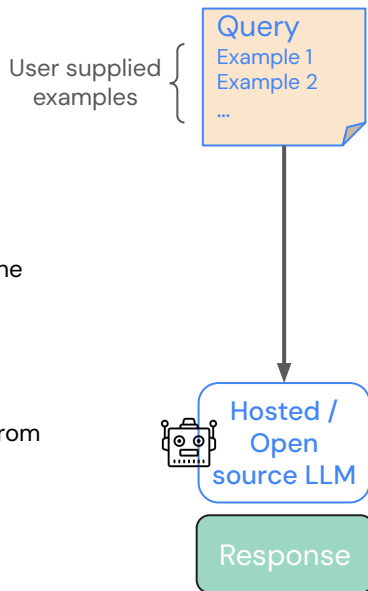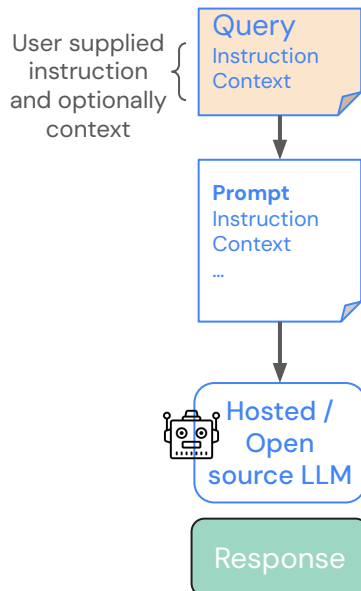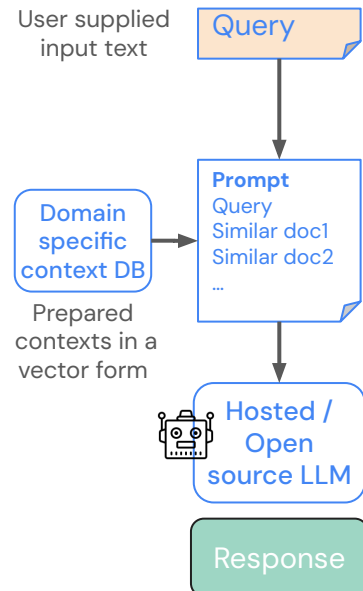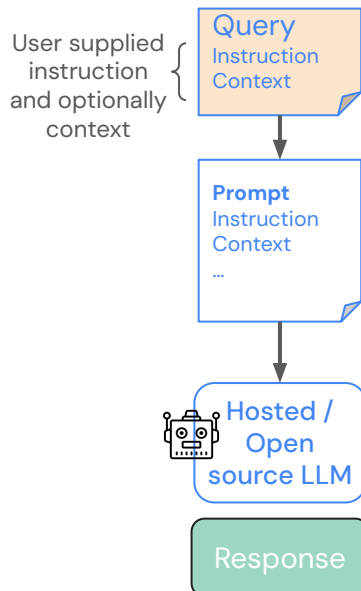This allows nesting or chaining LLMs, creating complex and dynamic interactions.

Hosted / Open source LLM

Response

**Note:** The patterns can be combined

# Prompt engineering patterns

"Add Context to the Query"

## Few shot prompting

User supplied examples {

**Query**
Example 1
Example 2
...

## Instruction following

User supplied instruction and optionally context

**Query**
Instruction
Context

**Prompt**
Instruction
Context
...

**Prompts can be:**

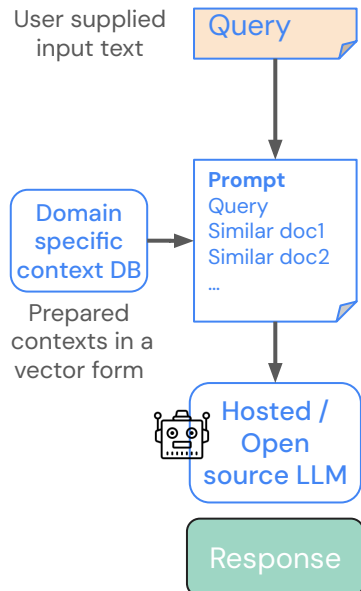Natural language sentences or questions. Code snippets or commands. Combinations of the above.

Emojis....basically any text!

Prompts can include outputs from other LLM queries.
This allows nesting or chaining LLMs, creating complex and dynamic interactions.

Hosted / Open source LLM

Response

Hosted / Open source LLM

Response

**Note:** The patterns can be combined

# Prompt engineering patterns

"Add Context to the Query"

## Few shot prompting

**Query**
Example 1
Example 2
...

User supplied examples

Hosted / Open source LLM

Response

## Instruction following

**Query**
Instruction
Context

User supplied instruction and optionally context

**Prompt**
Instruction
Context
...

Hosted / Open source LLM

Response

## Retrieval Augmented Generation

User supplied input text

**Query**

Domain specific context DB

Prepared contexts in a vector form

**Prompt**
Query
Similar doc1
Similar doc2
...

Hosted / Open source LLM

Response

**Prompts can be:**

Natural language sentences or questions. Code snippets or commands. Combinations of the above.

Emojis....basically any text!

Prompts can include outputs from other LLM queries.
This allows nesting or chaining LLMs, creating complex and dynamic interactions.

**Note:** The patterns can be combined

# Prompt engineering patterns

"Add Context to the Query"

**Few shot prompting**

User supplied examples

Query
Example 1
Example 2
…

Hosted / Open source LLM

Response

**Prompts can be:**

Natural language sentences or questions. Code snippets or commands. Combinations of the above.

Emojis….basically any text!

Prompts can include outputs from other LLM queries.
This allows nesting or chaining LLMs, creating complex and dynamic interactions.

**Instruction following**

User supplied instruction and optionally context

Query
Instruction
Context

**Prompt**
Instruction
Context
…

Hosted / Open source LLM

Response

**Retrieval Augmented Generation**

User supplied input text

Query

Domain specific context DB

Prepared contexts in a vector form

**Prompt**
Query
Similar doc1
Similar doc2
…

Hosted / Open source LLM

Response

**Note:** The patterns can be combined

# LLM Time

Policy Schedule / Validation Certificate

Darth Vader
54 Death Star, Nebula
London

(Please attach to policy document)

SINGLE CAR HIRE EXCESS INSURANCE POLICY

| | |
|---|---|
| Certification | No. BICEWCARTK/960 |
| Underwriter: | Newline Insurance Company Limited |
| Issue Date: | 16/07/2021 |
| Start Date: | 07/08/2021 Time: 11:00 AM |
| DropOff Date: | 14/08/2021 Time: 5:30 PM |
| Total Days | 8 Days |
| Area | Europe |
| Lead Driver First Name | Darth |
| Lead Driver Surname | Vader |
| Address | 54 Death Star, Nebula |
| | London |
| Email | darthvader@yahoo.com |
| Mobile Number | 7917077777 |

| First Name | Surname | Age |
|---|---|---|
| Darth | Vader | 245 |

Insurance Price: £17.86

**Insurance Policy pdf**

**txt**

**Vector Store**

Question + Context

**Prompt**

**LLM**

Embedded Question

Response

question

answer

50

# A Typical LLM Release

Multiple **sizes** (foundation/base model):

**small**

**base**

**large**

Size means memory required to load / train the model

# A Typical LLM Release

Multiple **sizes** (foundation/base model):

**large**

**base**

**small**

Size means memory required to load / train the model

Multiple **sequence lengths**:

512    4096    62000

Length you can learn from / use to generate text.

# A Typical LLM Release

Multiple **sizes** (foundation/base model):

**small**

**base**

**large**

Size means memory required to load / train the model

Multiple **sequence lengths:**

512　　4096　　62000

Length you can learn from / use to generate text.

Flavors/fine-tuned versions (**base**, **chat**, **instruct**):

I know what word comes next.

I know how to engage in conversation.

I know how to respond to instructions.

Models are trained on different instructions.

# What model are we gonna use?

# Application Time



**Insurance Policy pdf**

**txt**

**Vector Store**

Embedding

Question + Context

**Prompt**

Response

**LLM**

API

**question**

**answer**

# Chaining

# Chaining

# Q/A Bot



question

context

prompt

🦜🔗 LangChain

model

**+**

Answer
Handling

# Q/A Bot

# Agenda

- New Wave of Deep Learning
- Customisation Phases of GenAI
- Build a RAG application
- **Live demo**
- Fine Tuning Concepts
- 5 mins on Pretraining
- Do you really want to discuss cost?
- Summary – Call to Action

# RAG Architecture



**Insurance Policy pdf** → **txt** → **Vector Store** → Question + Context (**Prompt**) → **LLM**

Embedded Question

Response

question

answer

# KPIs Time

Foundational model as SaaS with Prompts

Retrieval Augmented knowledge (RAG)

Fine-tune foundational model on your data

Fully retrain foundational models (your own "GPT")
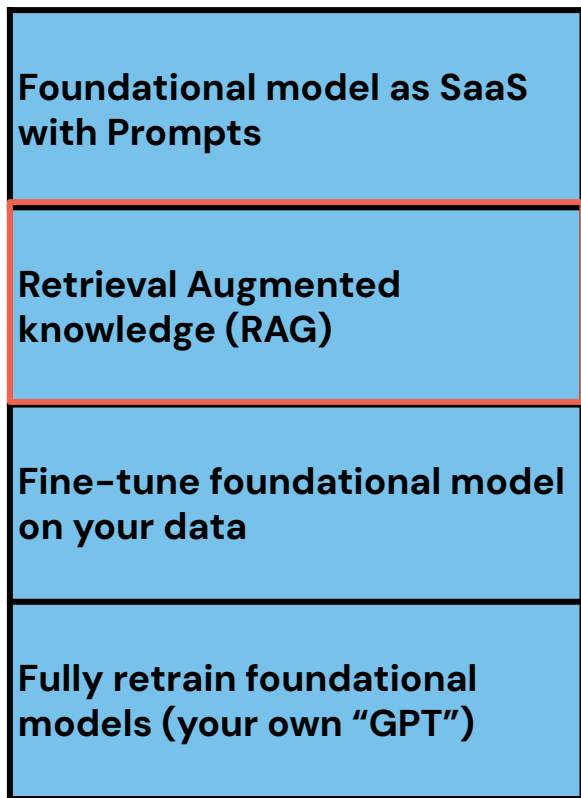
Training Cost

Data Size

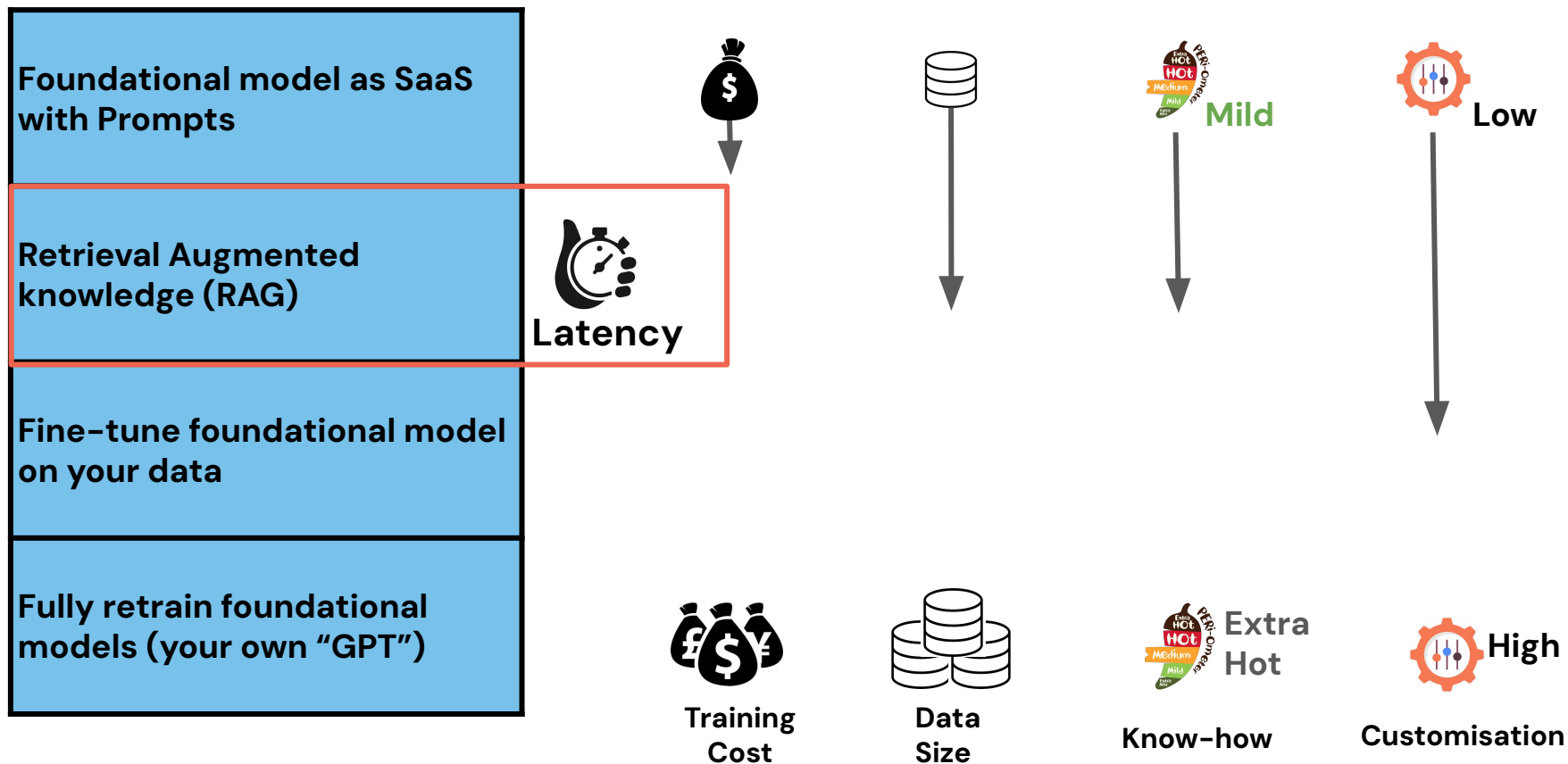Extra Hot
Know-how

High
Customisation

# Customisation Phases of GenAi

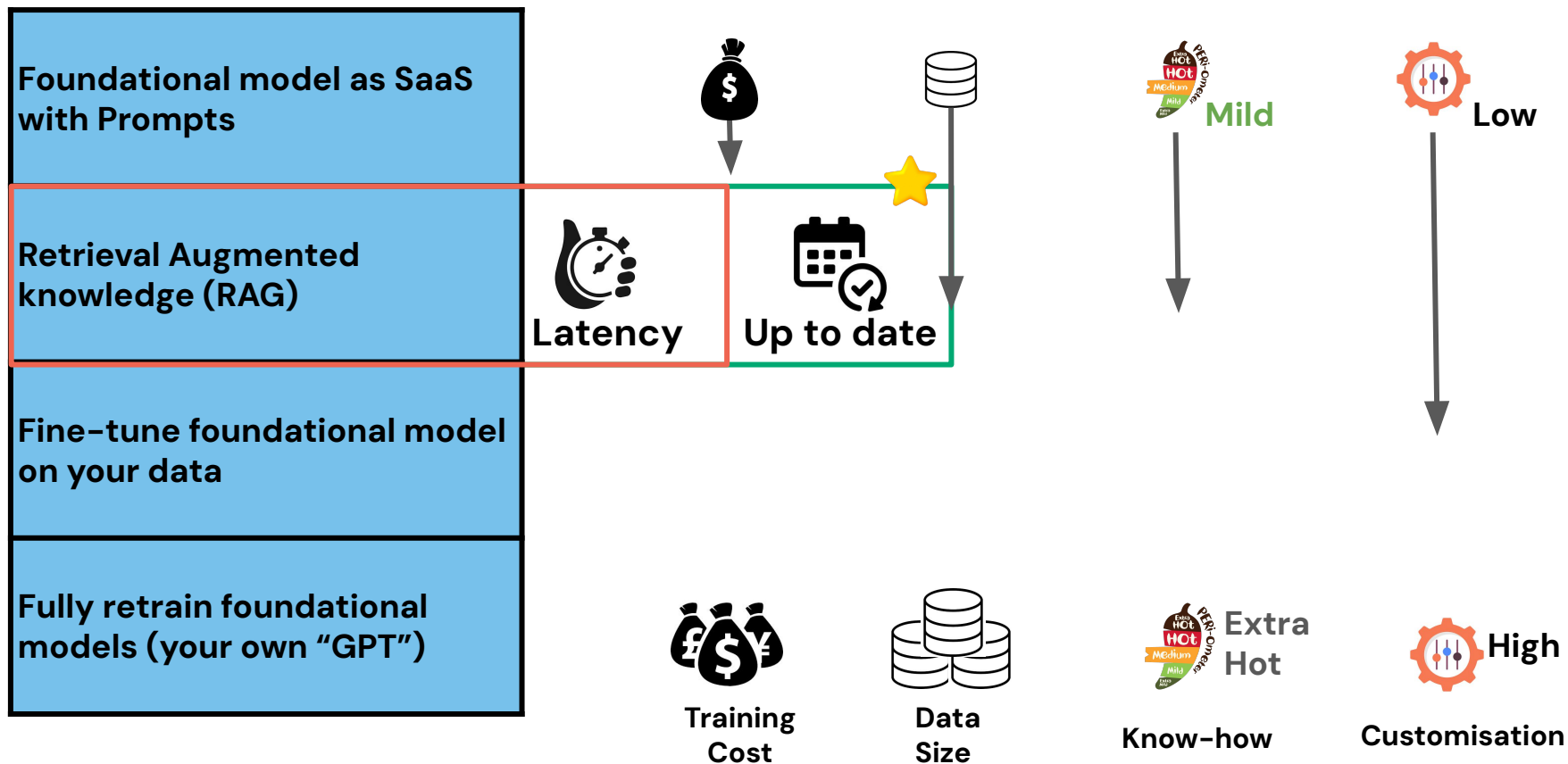| |
|---|
| **Foundational model as SaaS with Prompts** |
| **Retrieval Augmented knowledge (RAG)** |
| **Fine-tune foundational model on your data** |
| **Fully retrain foundational models (your own "GPT")** |

Mild

Low

Training Cost

Data Size

Extra Hot

High

Know-how

Customisation

# Customisation Phases of GenAi

| |
|---|
| **Foundational model as SaaS with Prompts** |
| **Retrieval Augmented knowledge (RAG)** |
| **Fine-tune foundational model on your data** |
| **Fully retrain foundational models (your own "GPT")** |

**Latency**

Mild

Low

**Training Cost**

**Data Size**

Extra Hot

High

**Know-how**

**Customisation**

# Customisation Phases of GenAi



| Foundational model as SaaS with Prompts |
| Retrieval Augmented knowledge (RAG) |
| Fine-tune foundational model on your data |
| Fully retrain foundational models (your own "GPT") |

Latency

Up to date

Mild

Low

Training Cost

Data Size

Know-how

Extra Hot

Customisation

High

# Phases of GenAi- RAG

| | |
|---|---|
| **Retrieval Augmented knowledge (RAG)** | - **Augment knowledge** of a GanAI model with **your own data**<br>- You can add filters to prompts **(avoid jailbreaking and hallucinations)**<br>- **Can control the model and version**<br>- Can **control ownership**<br>- - - - - - - - - - - - - - - - - - - - - - - - - - - -<br>- Still requires some **prompt engineering**<br>- You **don't control data inside** the model knowledge base<br>- It can **add latency** to your app |

# Final thoughts on RAG – Pros

- **Augment knowledge** of a GenAI model with your own data
- You can **add filters to prompts** (avoid jailbreaking and hallucinations)
- Can **control the model and version**
- Can **control ownership**

# Final thoughts on RAG – Cons

- It can get **expensive**
- You **don't control data inside** the model knowledge base
- It is not 100% clear **how the prompt affects the answer**
- **Domain specific Q/A** may **not work** well with RAG

# Agenda

- New Wave of Deep Learning
- Customisation Phases of GenAI
- Build a RAG application
- Live demo
- **Fine Tuning Concepts**
- 5 mins on Pretraining
- Do you really want to discuss cost?
- Summary – Call to Action

# Customisation Phases of GenAi – Fine Tuning

**Foundational model as SaaS with Prompts**

**Retrieval Augmented knowledge (RAG)**

**Fine-tune foundational model on your data**

**Fully retrain foundational models (your own "GPT")**

**Tune** a model on your data

# When should I fine tune models?

# Initial Process

Experimentation & Exploitation Strategy

# Initial Process

Experimentation & Exploitation Strategy

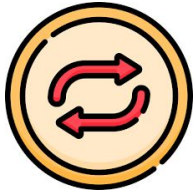# When should I fine tune?



**Repetition in the prompt – Token budget**

# When should I fine tune?

**Repetition in the prompt – Token budget**

**Promising few-shot**

# When should I fine tune?

**Repetition in the prompt – Token budget**
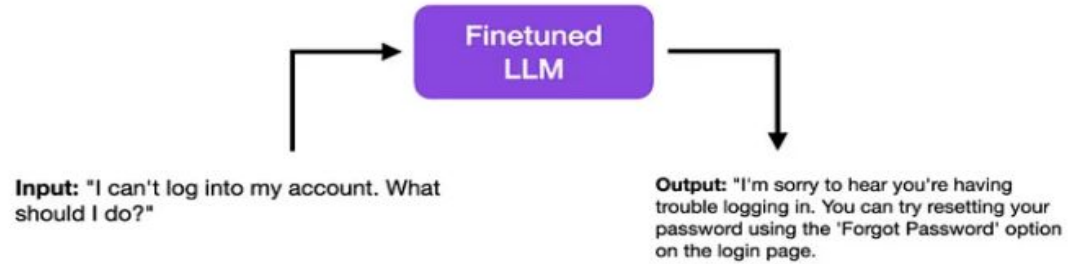
**Promising few-shot**

**Change the Behaviour**

# When should I fine tune?
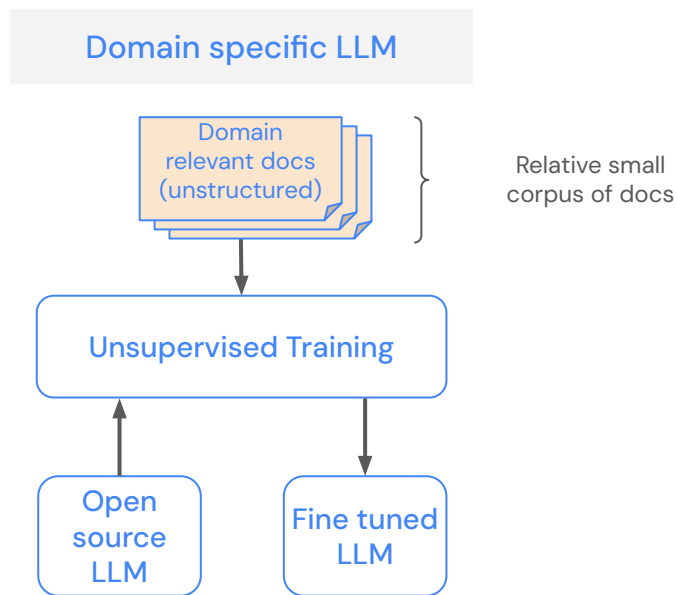
It's <u>NOT</u> for new concepts

# Fine Tuning - with an example



**Input:** "I can't log into my account. What should I do?"

**Pretrained LLM**

**Ouput:** "Try to reset your password using the 'Forgot Password' option."

# Fine Tuning - with an example



**Finetuned LLM**

**Input:** "I can't log into my account. What should I do?"

**Output:** "I'm sorry to hear you're having trouble logging in. You can try resetting your password using the 'Forgot Password' option on the login page.

# How are we fine tuning

# Fine-tuning Types - Domain Specific Tuning

"Adjust the model behavior"

# Fine-tuning Types - Instruction Tuning

"Adjust the model behavior"
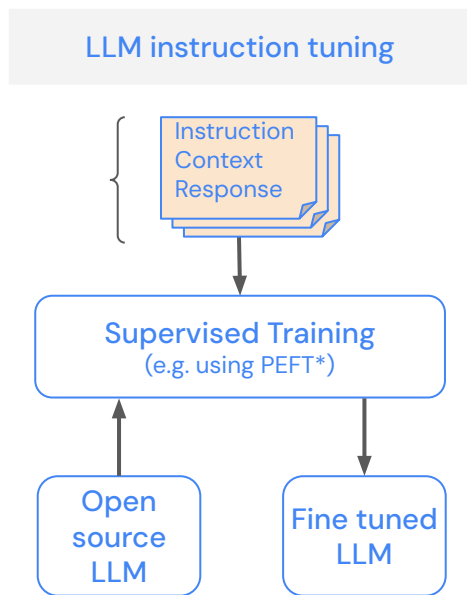
# Fine-tuning Types - Domain Specific Tuning

"Change the model behavior"

Domain specific LLM

Domain relevant docs (unstructured)

Relative small corpus of docs

Unsupervised Training (e.g. using PEFT*)

Open source LLM

Fine tuned LLM

- Fine tune on small corpus

# Fine-tuning Types - Instruction Tuning

"Adjust the model behavior"



**Notes:**

- LLM instruction tuning requires high quality labelled "instruction → response" data sets (increases effort & costs)

- Best results can be expected when combining both into two subsequent stages:
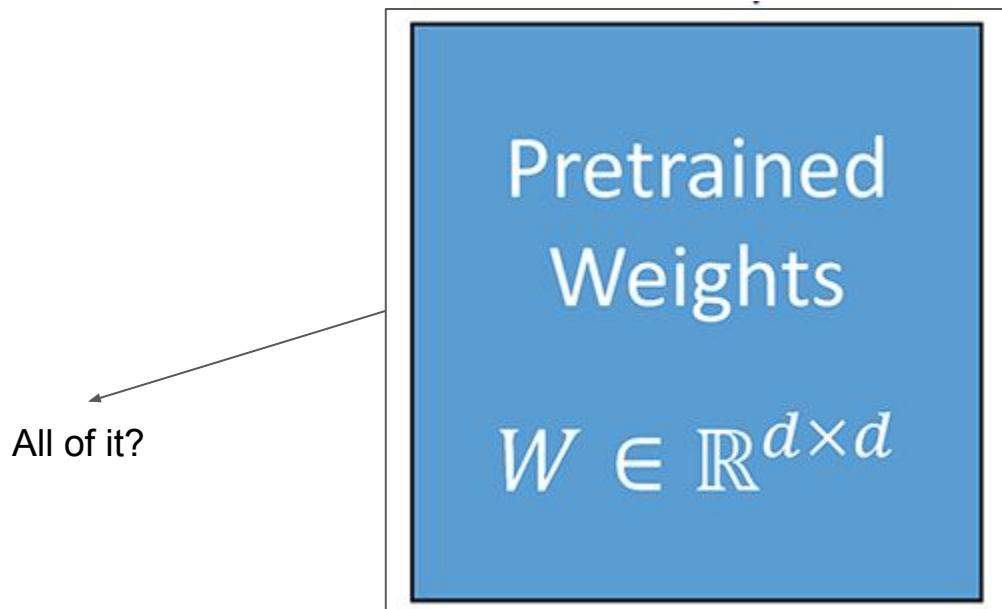
# What am I fine tuning?

# What am I fine tuning?

# What am I fine tuning?



Pretrained
Weights

$$W \in \mathbb{R}^{d \times d}$$

All of it?

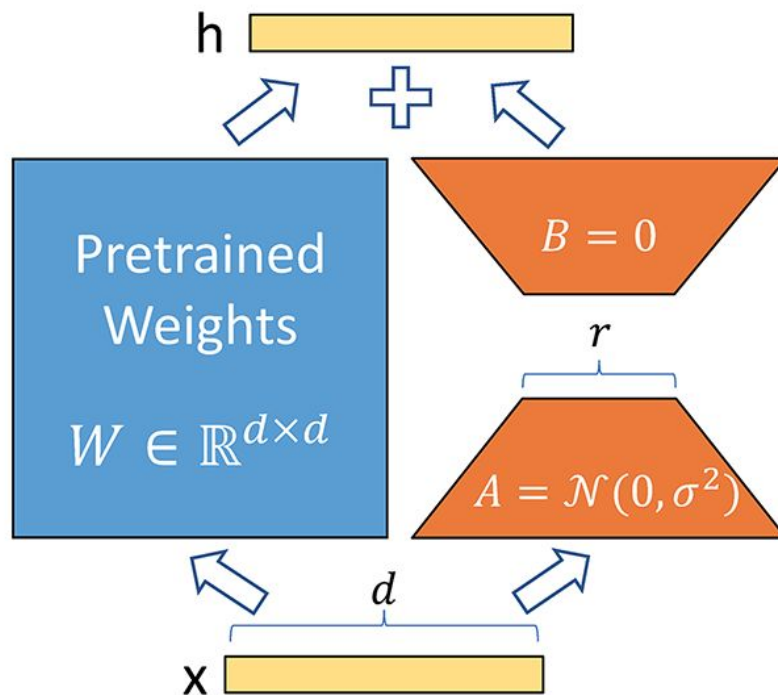# What am I fine tuning?



Pretrained Weights

$$W \in \mathbb{R}^{d \times d}$$
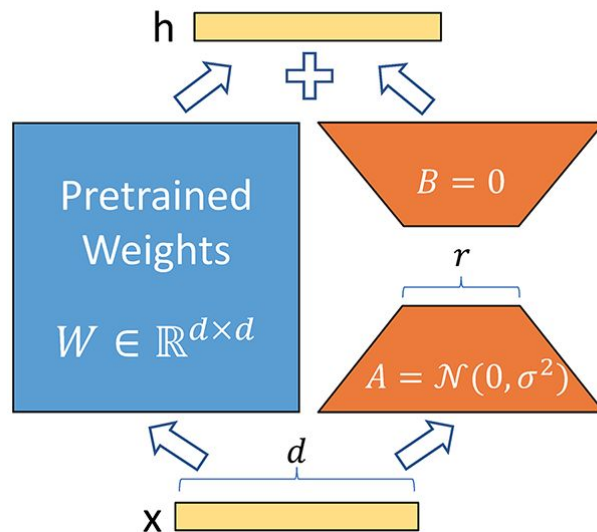
All of it?

1. Accelerate
2. Deepspeed

# What am I fine tuning?

# What am I fine tuning?
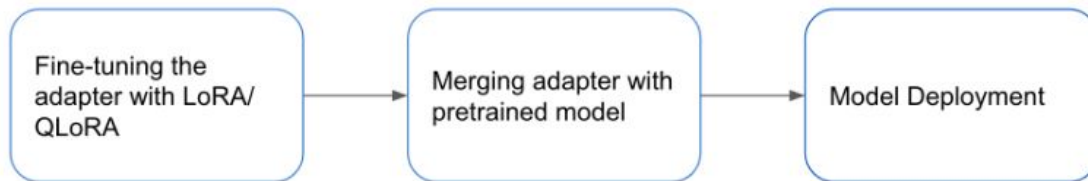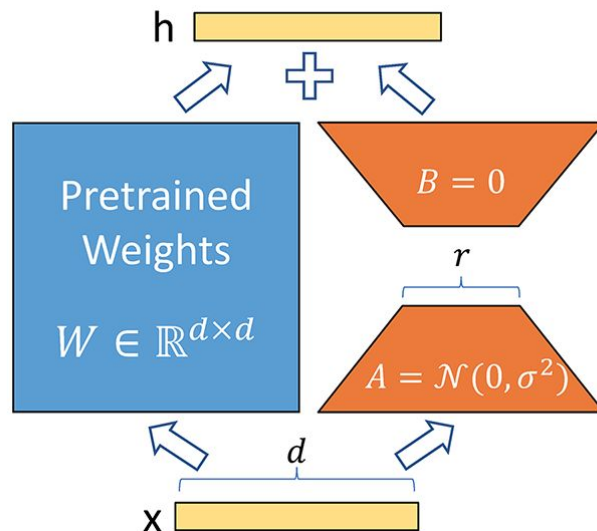
PEFT library → Parameter efficient Fine Tuning:

- **Lora** → Add Adapters with weights, which are the only parameters being fine-tuned, while freezing the rest
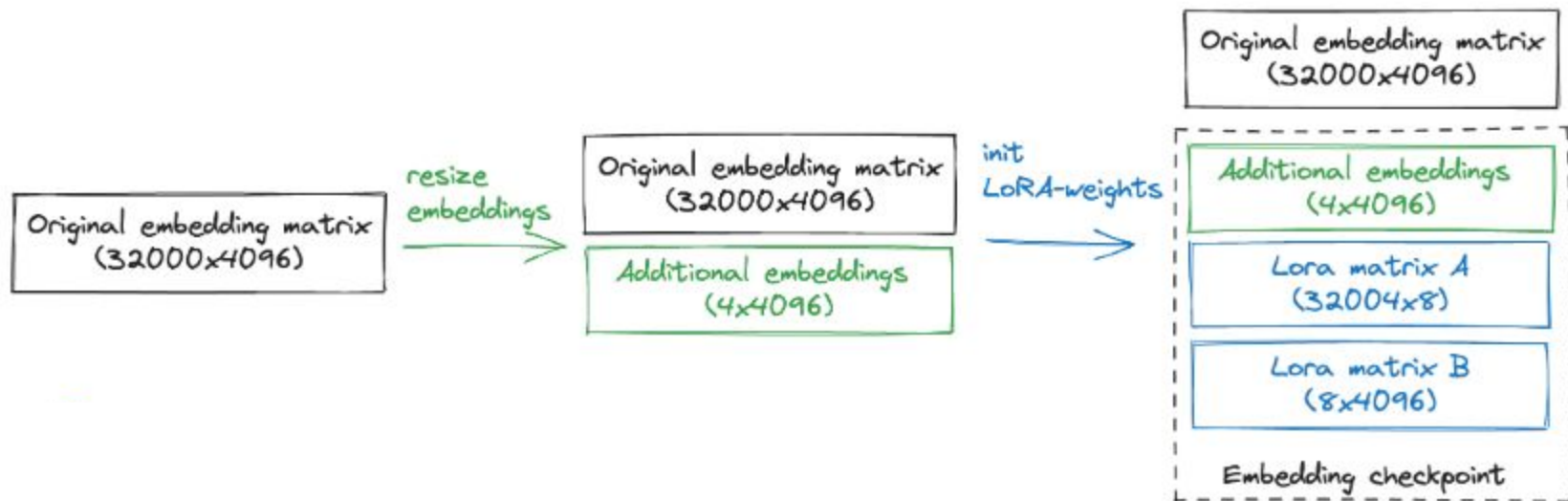- **Qlora** → As above but quantized version

# What am I fine tuning?

PEFT library → Parameter efficient Fine Tuning:

- Lora → Add Adapters with weights, which are the only parameters being fine-tuned, while freezing the rest
- Qlora → As above but quantized version

Original embedding matrix
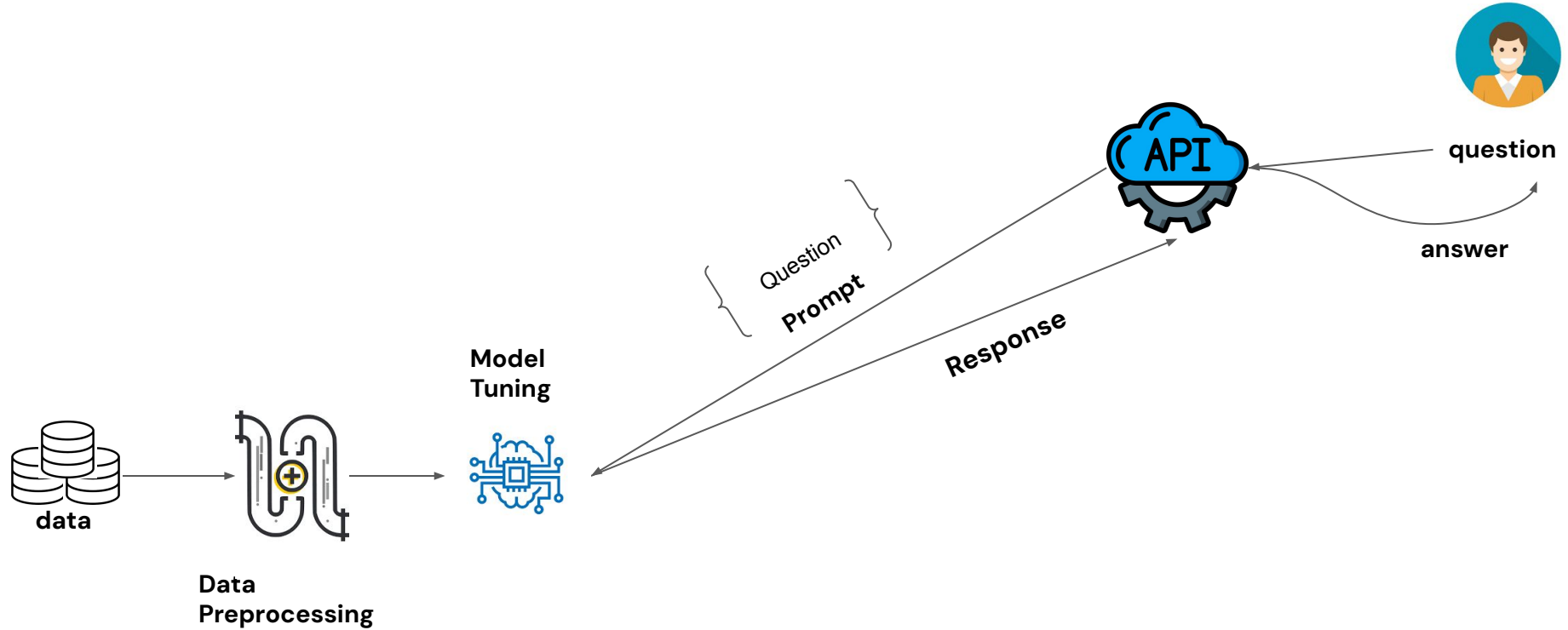(32000x4096)

resize
embeddings

Original embedding matrix
(32000x4096)

Additional embeddings
(4x4096)

init
LoRA-weights

Original embedding matrix
(32000x4096)

Additional embeddings
(4x4096)

Lora matrix A
(32004x8)

Lora matrix B
(8x4096)

Embedding checkpoint

# Checkpoint Sizes

[Ref](#)

| Number of trainable parameter / Checkpoint size | LoRA: q_proj and v_proj | LoRA: all layers | Full-parameter |
|---|---|---|---|
| 7B | 4194304 / 8MB | 20566080 / 41MB | 7B / 14GB |
| 13B | 6553600 / 13MB | 31887424 / 64MB | 13B / 26GB |
| 70B | 16384000 / 33MB | 104190016 / 201MB | 70B / 140GB |

# Final thoughts on Lora/Qlora

- The principal trade-off with LoRA is straightforward: you may give up some model quality, but you gain the ability to **serve many models more efficiently.**
- Cannot secure A100s? With LoRA you can still fine-tune models on **smaller GPUs** (reduced memory usage while training).
- Compared to regular checkpoints, LoRA checkpoints are significantly **smaller**, facilitating more **scalable serving, especially when managing multiple fine-tuned models.**

# Fine Tuning Architecture



data

Data
Preprocessing

Model
Tuning

Question
**Prompt**

**Response**

**question**

**answer**

# Fine Tuning Architecture With RAG



**Insurance Policy pdf**

txt

Embedding

API

question

answer

Response

Vector Store

Question + Context

Prompt

Tuned Model

More data

Data Preprocessing

# Kpis Time

| |
|---|
| **Foundational model as SaaS with Prompts** |
| **Retrieval Augmented knowledge (RAG)** |
| **Fine-tune foundational model on your data** |
| **Fully retrain foundational models (your own "GPT")** |

Training Cost

Data Size

Extra Hot

Know-how

High

Customisation

# Phases of GenAi

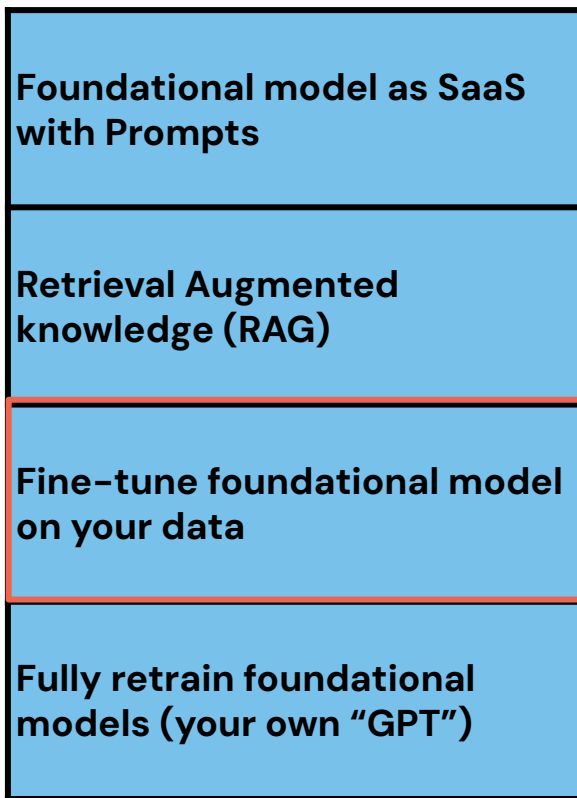| |
|---|
| **Foundational model as SaaS with Prompts** |
| **Retrieval Augmented knowledge (RAG)** |
| **Fine-tune foundational model on your data** |
| **Fully retrain foundational models (your own "GPT")** |

Mild

Low

Hot

**Training Cost**

**Data Size**

**Know-how**

High

**Customisation**

# Phases of GenAi– Fine Tuning

| Fine–tune foundational model on your data | - Can **update certain "parts"** of the model |
|---|---|
| | - Can **win with a smaller model** |
| | - Can still **add RAG and** add filters to prompts **(avoid jailbreaking and hallucinations)** |
| | - **Can control the model and version** |
| | - Can **control ownership** |
| | - Still requires some **prompt engineering** |
| | - You **don't control data inside** the model knowledge base |
| | - **No guarantee this can improve quality** |
| | - Requires **computational resources and technical skills** |

# Final thoughts on Fine Tuning – Prons



- Trained on domain specific knowledge so more accurate responses (may, may not)
- You can **lock down the version of the model and IP**
- Can **still add RAG** and add filters to prompts (avoid jailbreaking and hallucinations)
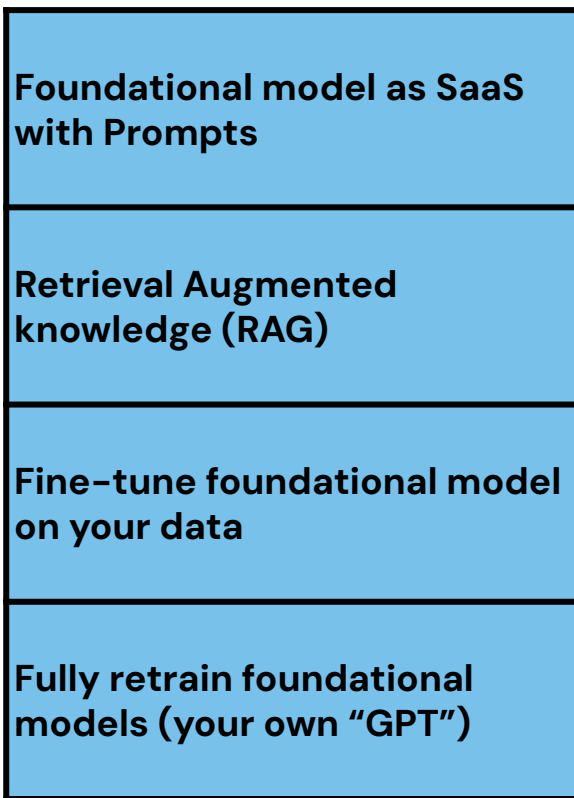
# Final thoughts on Fine Tuning – Cons



- You need to **gather the data** and make sure **they are of good quality**
- Fine-tuning typically results in creating a **niche model for a niche use-case**
- Model management and infrastructure, and serving
- Original pretraining data **may dominate**
- Can't create new capability, just bring domain specific knowledge to the model
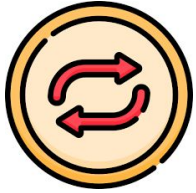- No **guarantee** this can improve quality

# Agenda

- New Wave of Deep Learning
- Phases of GenAI
- Build a quick RAG application with your own data
- Live demo
- Fine Tuning Concepts
- Live demo
- **5 mins on Pretraining**
- Do you really want to discuss cost?
- Summary – Call to Action

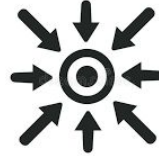# Customisation Phases of GenAi – Pretraining

**Foundational model as SaaS with Prompts**

**Retrieval Augmented knowledge (RAG)**

**Fine-tune foundational model on your data**

**Fully retrain foundational models (your own "GPT")**

**Create** a model on your data

# Why Pretraining:

**Full customisation**

**Full IP**

**Competitive Advantage**

# Why pretrain models?

**Control**
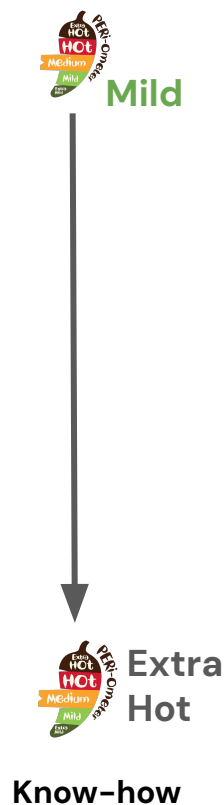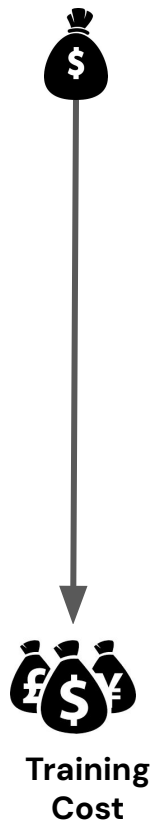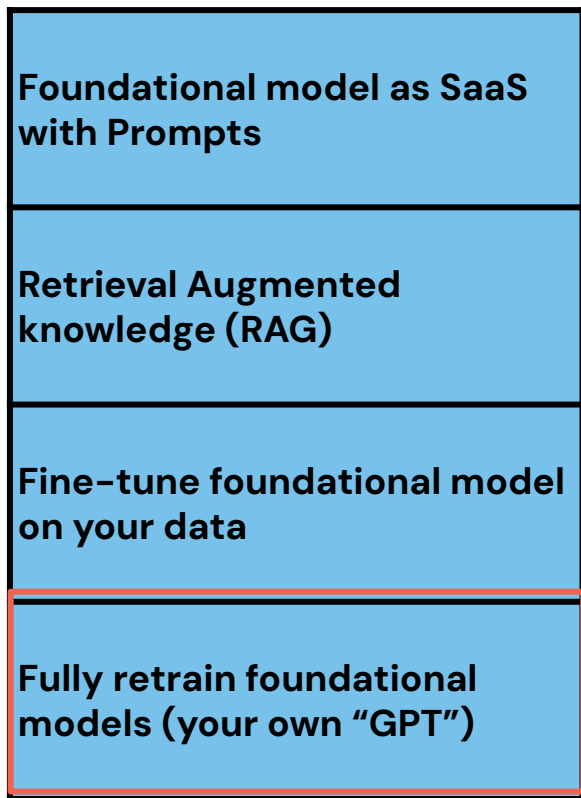
To Train Generative AI models we need…..

GPUs

**We need a lot of GPUS to train <u>your own</u> Generative AI models**

# Kpis

| |
|---|
| **Foundational model as SaaS with Prompts** |
| **Retrieval Augmented knowledge (RAG)** |
| **Fine-tune foundational model on your data** |
| **Fully retrain foundational models (your own "GPT")** |

Mild

Low

Extra Hot

High

**Training Cost**

**Data Size**

**Know-how**

**Customisation**

# Phases of GenAi- RAG

| | |
|---|---|
| **Retrieval Augmented knowledge (RAG)** | – **Owner** of your mini GPT style<br>– The way to go for a very **particular use cases or small models**<br><br>– Requires **resources** both technical and computational<br>– Requires a **lot of data** or labels (100K+)<br>– If you don't know what you are doing, it will **never converge** |

# Final thoughts on Fine Tuning – Prons



- Owner of your **mini GPT style**

- You **control data** inside the model knowledge base

# Final thoughts on Fine Tuning – Cons

- Requires **resources both technical and computational**

- Requires **a lot of data**

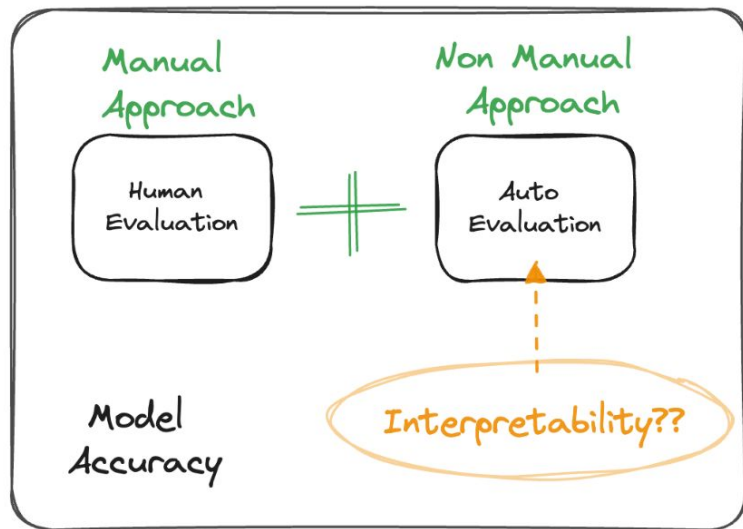- If you don't know what you are doing, **it will never converge**

# Agenda

- New Wave of Deep Learning
- Phases of GenAI
- Build a quick RAG application with your own data
- Live demo
- Fine Tuning Concepts
- Live demo
- 5 mins on Pretraining
- **Do you really want to discuss cost?**
- Summary – Call to Action

# LLM Projects: Cost

- **Instruction Fine Tune:**
  - Start from 1–10k training examples
  - PEFT

- **Continue Pre–Training**
  - Starts from around 100m–1bn tokens
  - PEFT

- **Pre–Train LLM from scratch:**
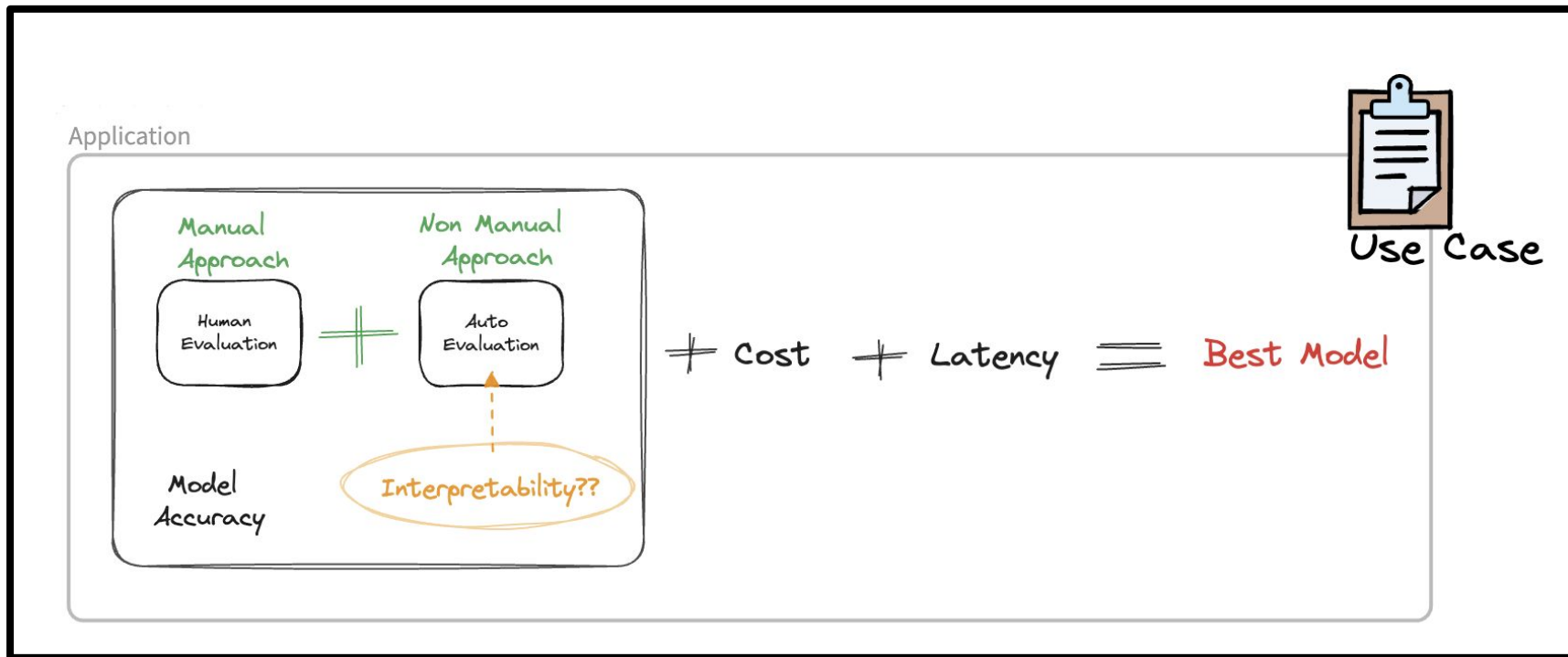  - Requires carefully crafted and very huge(1T) training datasets

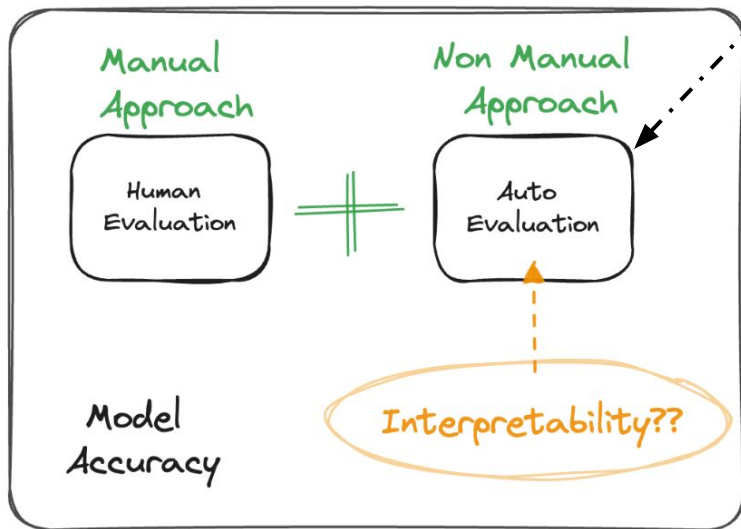| Complexity | Cost | Training Loop |
|---|---|---|
| Distributed GPU setup required High | Depending on model size Medium to High | Occasional |
| Distributed GPU setup required High | Depending on model size Medium to High | Occasional |
| Distributed GPU setup required High | 100K – 2.5 Mil $ Very high | Rare |

# Your Best LLM

# Your Best LLM

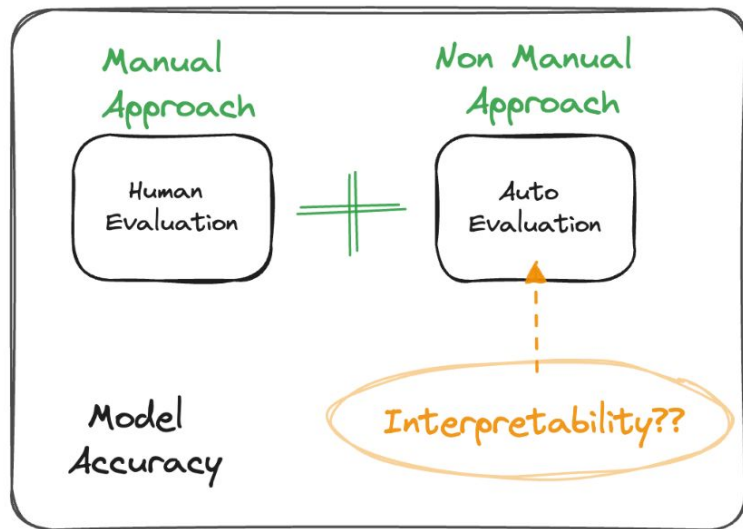# Evaluate LLMs – Auto Evaluation

Translation (BLEU),
Summarization (ROUGE)
Q&A (F1, toxicity, perplexity, exact match)
Document retrieval (NDCG, MMR)

Manual Approach

Non Manual Approach

Human Evaluation + Auto Evaluation

Model Accuracy

Interpretability??

+ Cost + Latency = Best Model

# Your Best LLM

# Agenda

- New Wave of Deep Learning
- Customisation Phases of GenAI
- Build a quick RAG application with your own data
- Live demo
- Fine Tuning Concepts
- Live demo
- 5 mins on Pretraining
- Do you really want to discuss cost?
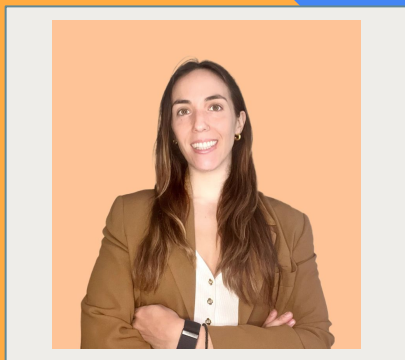- **Summary – Call to Action**

# Summary

- New wave of LLMs

- Customisation Phases of LLMs

  - Prompt Engineering

  - RAG

  - Fine Tuning

  - Pretraining

- Aspects of cost

# An Overview of Common LLMs

| Use case | Quality-optimized | Balanced | Speed-optimized | Notes |
|---|---|---|---|---|
| Text generation following instructions | Mixtral-8x7B-Instruct-v0.1<br><br>MPT-30B-Instruct †<br><br>Llama-2-70b-chat-hf | Mistral-7B-Instruct-v0.2<br><br>MPT-7B-Instruct †<br><br>MPT-7b-8k-instruct<br><br>Llama-2-7b-chat-hf<br><br>Llama-2-13b-chat-hf | phi-2 | † Supervised fine-tuning using databricks-dolly-15k dataset |
| Text embeddings (English only) | e5-mistral-7b-instruct (7B) | Bge-large-en-v1.5 (0.3B)<br><br>e5-large-v2 (0.3B) | bge-base-en-v1.5 (0.4B)<br>bge-small-en-v1.5 (0.1B)<br>e5-base-v2 (0.1B) | |
| Transcription (speech to text) | | whisper-large-v3 (1.6B) | distil-large-v2 (0.7B) | |
| Image generation | | stable-diffusion-xl | | |
| Code generation | CodeLlama-70b-hf<br><br>CodeLlama-70b-Instruct-hf<br><br>CodeLlama-70b-Python-hf (Python optimized)<br><br>CodeLlama-34b-hf<br><br>CodeLlama-34b-Instruct-hf<br><br>CodeLlama-34b-Python-hf (Python optimized) | CodeLlama-13b-hf<br><br>CodeLlama-13b-Instruct-hf<br><br>CodeLlama-13b-Python-hf (Python optimized)<br><br>CodeLlama-7b-hf<br><br>CodeLlama-7b-Instruct-hf<br><br>CodeLlama-7b-Python-hf (Python optimized) | | Code LLMs usually need fine-tuning to follow instructions and work on application-specific code |

# Thank you so much!

Maria Zervou

Sr. Specialist Solutions Architect

Code & Slides