









Welcome to the AI Capability Hands-on tailored for SMEs.

This interactive session is tailored to equip both technical and non-technical individuals with a thorough grasp of GenAI, its practical uses, and essential guidelines, particularly focusing on containment and alignment. As GenAI becomes increasingly integrated into our daily lives, establishing robust guardrails will be pivotal for responsible deployment.

During the program, attendees will explore the complexities of various models and frameworks, grapple with the ethical considerations of AI, and gain hands-on experience in developing their own GenAI applications, agents, and chatbots, among other subjects.





Understanding AI and Its Foundations

First part focuses on introducing participants to AI and its foundational concepts. The sessions are designed to be interactive, with hands-on exercises complementing theoretical knowledge.

The highlights include:

- o Demystifying Models, Weights, and Biases
- o The differences between tunning vs Prompting
- o Introduction to AI frameworks with a spotlight on Langchain
- Exploring no-code AI solutions with LangFlow and Flows





Implementing, Testing, and Safeguarding AI

Part two shifts the focus towards implementing, testing, and safeguarding AI solutions. Participants will gain insights into testing AI models, addressing safety concerns, and understanding the ethical considerations in AI.

The highlights include:

- Testing and evaluating models using tools like LangSmith and LangForge
- o Introduction to GuardRails for AI to ensure business safety
- o Practical hands-on session to build a simple AI project



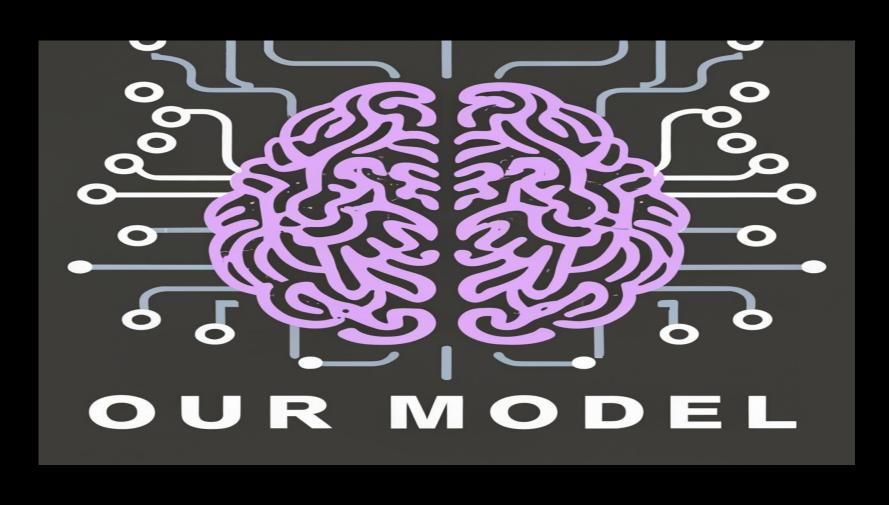


Key Takeaways at the end of this interactive session, participants will

- Have a clear understanding of AI, models, weights, and biases.
- o Be familiar with various AI frameworks and tools.
- o Gain hands-on experience in building AI applications.
- Understand the importance of safety and ethics in AI.
- Be equipped with the knowledge to scale AI solutions for their businesses.



What is a Model







Your Objectives & Key Results

Inference <u>Call</u>

Like finally baking the perfect cookies and sharing them

Prompting

Like asking Can you make the cookies chocolate flavoured

Tuning

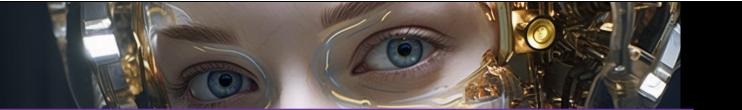
Like baking and tasting cookies to make them perfect.

Weight & Biases

Like the amounts of sugar and flour. Adjust them to make cookies taste better.

Model

Think of this like the recipe for cookies. It tells us the ingredients and steps.





Tunning Vs Prompting





Tuning

Modular Fine-Tuning of Pre-existing Models For domain-specific tasks, Fine Tuning of existing models is vital. It customizes pre-trained models for specific tasks using domain data. However, this can be costly, as it modifies millions of parameters, needing extensive data and infrastructure. Full fine-tuning risks catastrophic forgetting and results in large, non-modular models. Solutions like Parameter Efficient Fine Tuning (PEFT) offer modular fine-tuning that conserves resources and is cost-efficient.

Tunning:(PEFT-LoRA, QLoRA)

Llama 2 GPT-3-Turbo PALM 2 Your Objectives & Key Results GenAI Frameworks & Applications Your Models LLM

Prompting

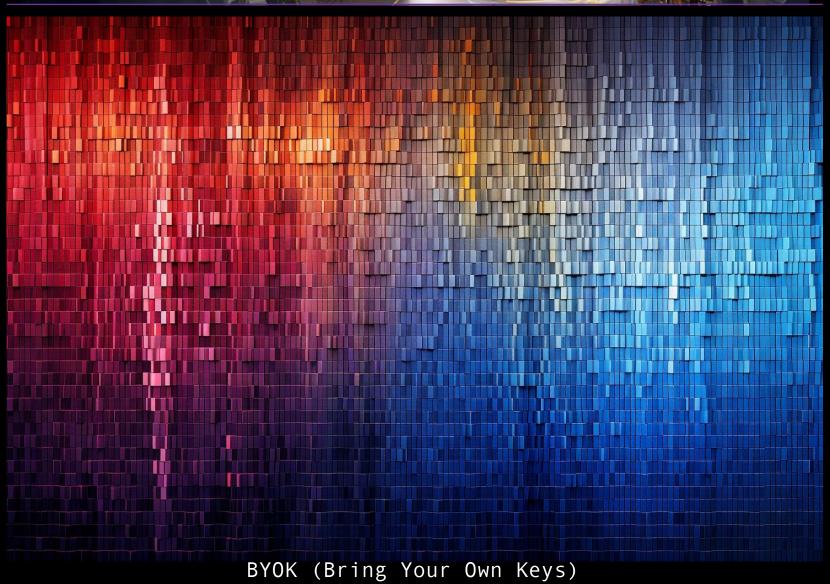
Prompt Engineering Preexisting Models Fundamentally, obtaining anticipated results from Large Language Models requires meticulous engineering of prompts. This method entails designing appropriate prompts and inputs to draw out the specific reactions from the model. Engineering prompts is a crucial approach for numerous applications, particularly when general feedback is adequate.

Prompting: (MRKL, REAct, COT, RAG, Zero Shot)

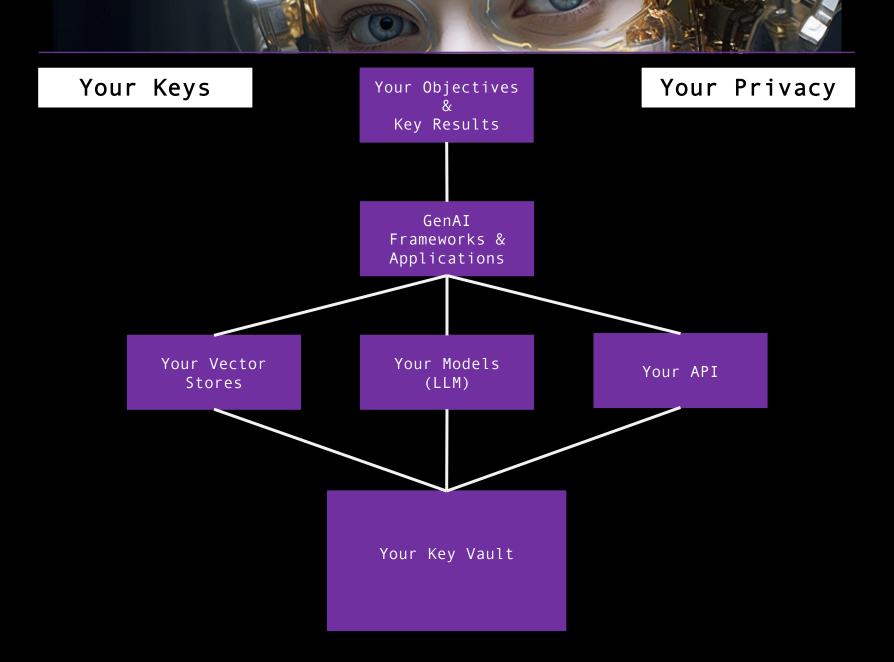
Llama 2 GPT-3-Turbo GPT-4 PALM 2 Claude 2







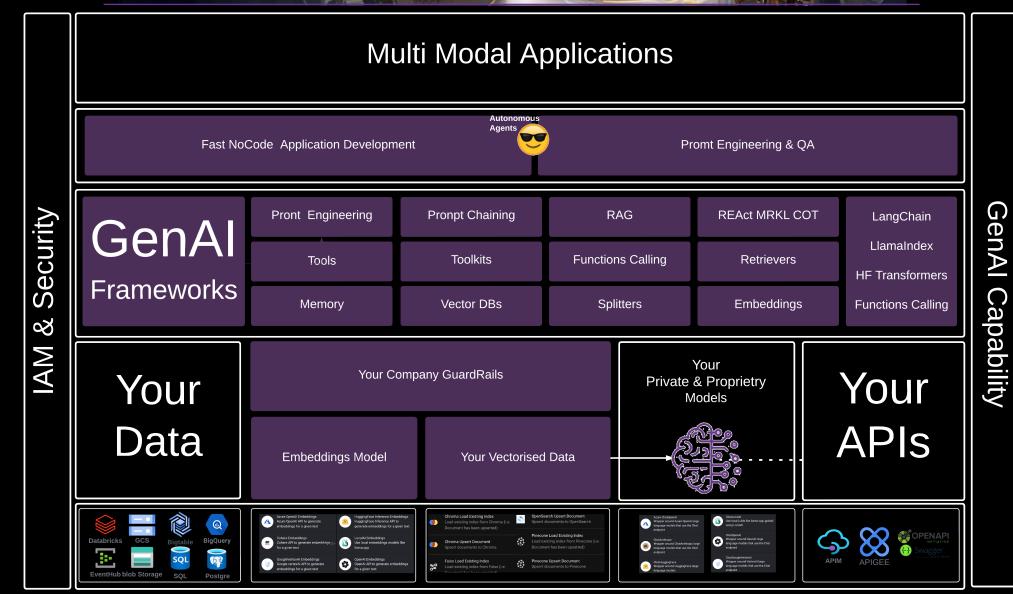


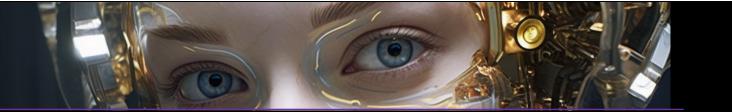












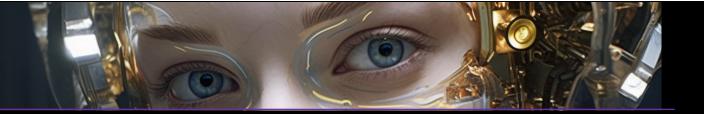


Glossary Terms

```
Attention Model architecture Completion Low rank adaptation (LoRA) Fine tuning
Artificial Intelligence (AI)
Reinforcement Learning from Human Feedback (RLHF) One-shot / Few-shot
        Parameters Transformer Chain-of-thought
Generative pretrained transformers (GPT
Generative AI Embeddings

Multi-modal Multi-modal
Large language model (LLM) Trainir
Agents

Agents
Agents
       Neural network Token ChatGPT Foundational model Alignment
```





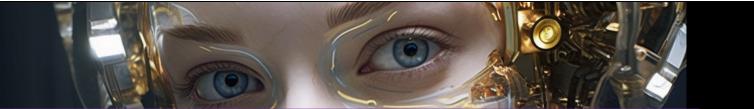
Term	Description
Neural network	Network modelled on the brain
Parameters	Weights that control neural network calculations
Model architecture	Components of a complex AI model
Training	Improving model performance on data
Generative AI	Models that generate text/images from prompts
Generative pretrained transformers (GPT)	Popular large language model
ChatGPT	Conversational version of GPT
Large language model (LLM)	AI model that handles language
Transformer	Popular neural network architecture
Token	Encodes text numerically for models
Embeddings	Represent words/text semantically
Attention	Allows models to understand context



Term	Description	
Alignment	Steers models towards ethical output	
Foundational model	Broadly trained model	
Fine tuning	Tailoring model to specific tasks	
RLHF	Reinforcement Learning from Human Feedback to improve models	
Low rank adaptation (LoRA)	Efficient fine-tuning method	
Multi-modal	Handle mixed text/image input	
Prompt	Text input to models	
Completion	Text output from models	
Hallucination	Fictional/incorrect output	
One-shot / Few-shot	Types of prompting	
System prompt	Defines model characteristics	
Prompt engineering	Developing effective prompts	

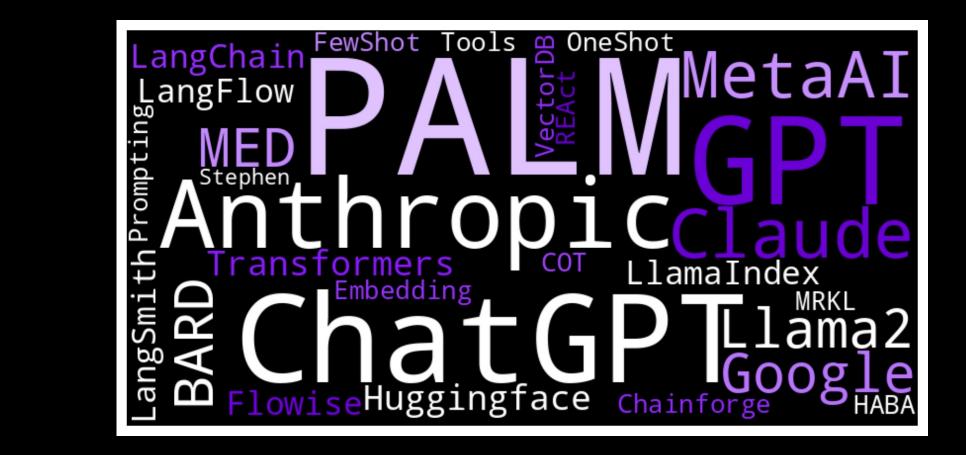


Term	Description	
Prompt injection	Subverting models via input	
Chain-of-thought	Improves reasoning via breakdown	
REAct	Read Evaluate and React	
MRKL	Modular Reasoning, Knowledge & Language	
Agents	Versatile AI tools	
Plugins / tools	Expand capabilities via APIs	
Retrieval Augmented Generation (RAG)	Supplementing with searches	





Appendix





Term	Description	Links
Nvidia Guardrails	Stop AI systems getting out of control	
<u>Langchain:</u>	 Chaining large language models together, allowing them to converse 	https://www.langchain.com
Flowise	Flowise: Workspace for building AI flows visually	https://flowiseai.com
Chainlit	Service for running computations on GPT models	https://docs.chainlit.io/o verview
LangFlow	Framework for creating AI pipelines with Python	https://www.langflow.org
LangSmith	Platform for searching AI models and chaining them together	https://smith.langchain.co m
Retrieval Augmented Generation (RAG)	Supplementing with searches	
LLM's	(GPT-4, Anthropic, PALM 2, Llama 2)	<pre>https://openai.com https://www.anthropic.com https://cloud.google.com/v ertex-ai https://ai.meta.com</pre>