

# INTRODUCTION TO LARGE LANGUAGE MODELS (LLMs)

Yogesh Haribhau Kulkarni



# Outline

1 AI INTRO

2 GENAI INTRO

3 PROMPTENGG DEMO

4 CONCLUSIONS

5 REFERENCES

# Introduction to Artificial Intelligence

YHK

“Houston, we have a problem!!”



50 Years Ago: “Houston, We’ve Had a Problem” – John Uri

YHK

## Whats the Problem?

- ▶ Along with some softer words like “disruption”, “passionate”, “excited” ...
- ▶ If you don't have word “innovation” in your talk/speech/conversation it's BIG problem.
- ▶ Irrespective of fields. You can be Corporate, Political, Social, etc.

And there is an addition of one more word, which is a must in every talk...and that is?

## The Problem

Every company is claiming to be working in AI-ML

- ▶ Is it really so?
- ▶ What exactly is AI (ML)?
- ▶ What is not AI?

Or is it just a plain BIG hype?



# What is the Core Idea?

YHK

## What's the core idea?

- ▶ behind problem solving?
- ▶ behind writing software algorithms?
- ▶ solving research problems?



# Desire

- ▶ To find a “function”
- ▶ To find a relation
- ▶ To find a transformation
- ▶ To build a model
- ▶ From given inputs to desired outputs.  
That's it.



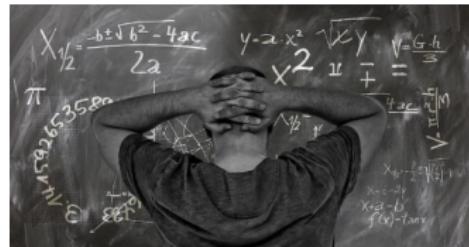
## Functions

- ▶ Some functions are straight forward
- ▶ *"In summer, ice-cream sale goes up"*
- ▶ Cause and effect
- ▶ Relation (function, Mathematical model) is found out
- ▶ Here, simple rule based programming suffices



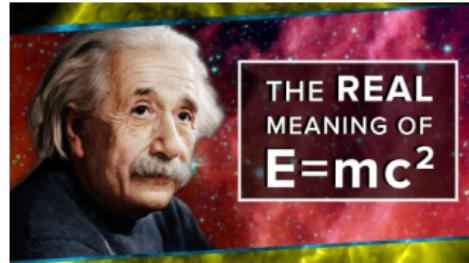
# Functions

- ▶ But some functions are complex
- ▶ *"More you put efforts, your business flourishes."*
- ▶ Cause and effect again, but the relation is far too complex
- ▶ Too many variables
- ▶ Here, simple rule based programming not humanly possible.
- ▶ Lots of research needed to come up with equations.



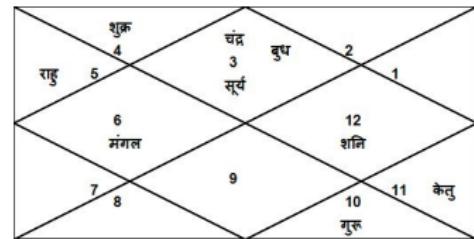
# Functions

- ▶  $E = mc^2$
- ▶ What's this? a function?
- ▶ Input variable(s)?
- ▶ Output variable(s)?
- ▶ Parameters?
- ▶ How's the relation? linear?



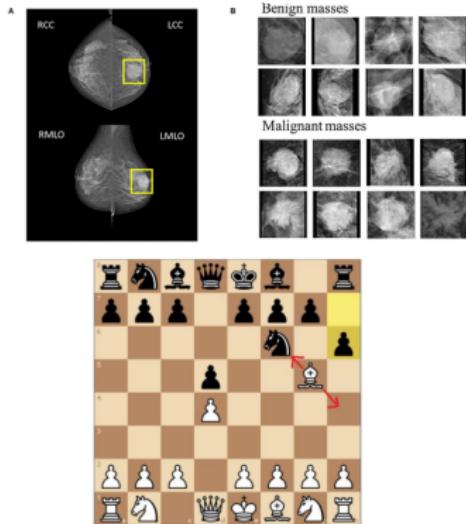
## Controversial Example

- ▶ Even astrology is a model, based on the past cases.
- ▶ Could claim empirical evidence.
- ▶ Given this planetary position, it predicts.
- ▶ Represented by "Horoscope"
- ▶ Got weights for each planets (real or fictitious)
- ▶ Reliable??



## Functions

- ▶ But most real-life functions are not deterministic
- ▶ Some are probabilistic, some non-linear.
- ▶ “*Detecting if the tumor is benign or malignant*”
- ▶ “*At any state in the game of chess, what's the next move?*”



## Chess: next move?

- ▶ Needs extreme expertise
- ▶ Needs “intelligence”
- ▶ How do you get that?
  - ▶ Built by lots of training.
  - ▶ By studying lots of past games.
- ▶ This is how Humans build intelligence



# Intelligence

- ▶ Can machine (software/program) also do the same?
- ▶ Can it play chess?
- ▶ Can it build intelligence?
- ▶ By looking at past experiences (data),
- ▶ Training Data: games played, moves used, etc.

Yes, it can!! That's Artificial Intelligence.



# What is Artificial Intelligence?

YHK

## My definition

“If machines (or computer programs) start doing some/all of these “intelligent” tasks, then that’s Artificial Intelligence”

## Intelligence: the differentiation

- ▶ Ability to think various domains
- ▶ Ability produce something new
- ▶ Ability to detect the unseen
- ▶ Ability to enhance knowledge (rules, patterns)



All these, AI has started doing. The AI era has arrived!!

## Everyday usage

Artificial intelligence seems to have become ubiquitous.

- ▶ Replying to our emails on Gmail
- ▶ Learning how to drive our cars,
- ▶ Sorting our holiday photos.
- ▶ etc.



Too good to be true, isn't it, sort of Magical !!

## But then ...

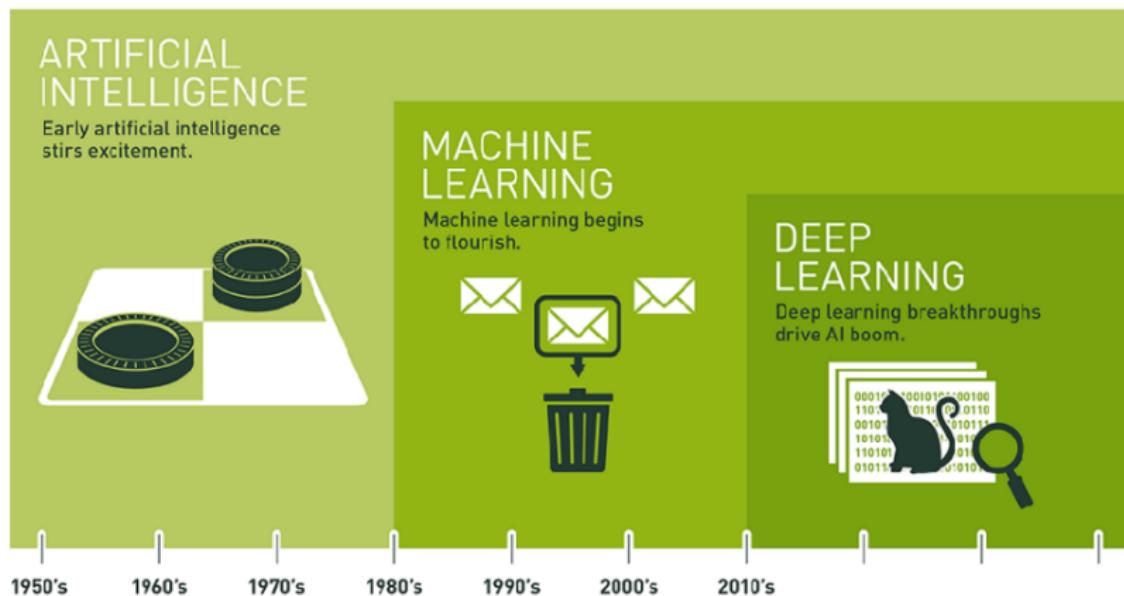
- ▶ When its too good, you start suspecting
- ▶ Is it for real!!
- ▶ How can such thing happen?
- ▶ How far will it go?



The next thing you know, people are worrying about exactly how and when AI is going to doom humanity.

# AI, ML, DL ... Same?

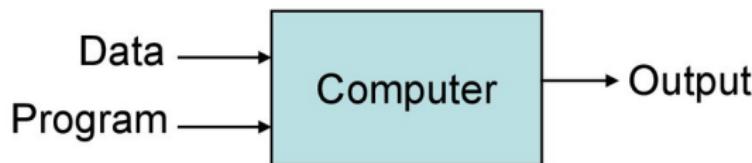
Or Relationship between them ?



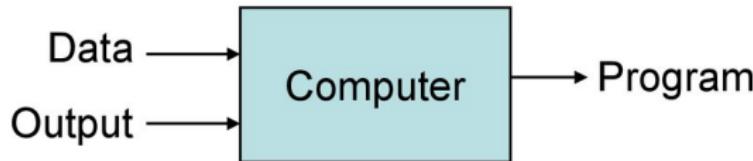
(Ref: <https://blogs.nvidia.com/blog/2016/07/29/whats-difference-artificial-intelligence-machine-learning-deep-learning-ai/>)

# Traditional vs. Machine Learning?

## Traditional Programming



## Machine Learning



# Why Machine/Deep Learning?

- ▶ Problems with High Dimensionality
- ▶ Hard/Expensive to program manually
- ▶ Techniques to model 'ANY' function given 'ENOUGH' data.
- ▶ Job \$\$\$

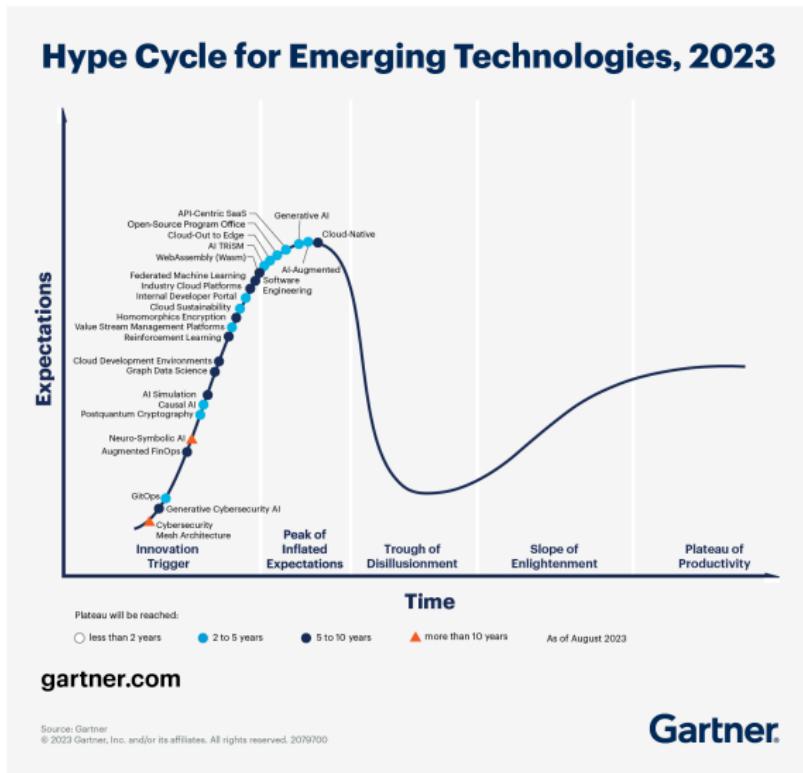


## Why now?

- ▶ Flood of data (Internet, IoT)
- ▶ Increasing computational power
- ▶ Easy/free availability of algorithms
- ▶ Increasing support from industries



# Gartner Hype Cycle Emerging Technologies 2023



# Is AI a threat?



## Is AI a threat?

If you believe in what Elon Musk says, then YES.



*Elon Musk recently commented on Twitter that artificial intelligence (AI) is more dangerous than North Korea*

(Ref: What is Artificial Intelligence — Artificial Intelligence Tutorial For Beginners — Edureka)

# Is AI a threat?

If you believe in these movies, then YES.



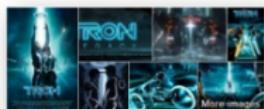
The Terminator



I, Robot



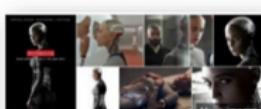
The Matrix



Tron: Legacy



War Games



Ex Machina

Well, AI based War robots are not impossible anymore.

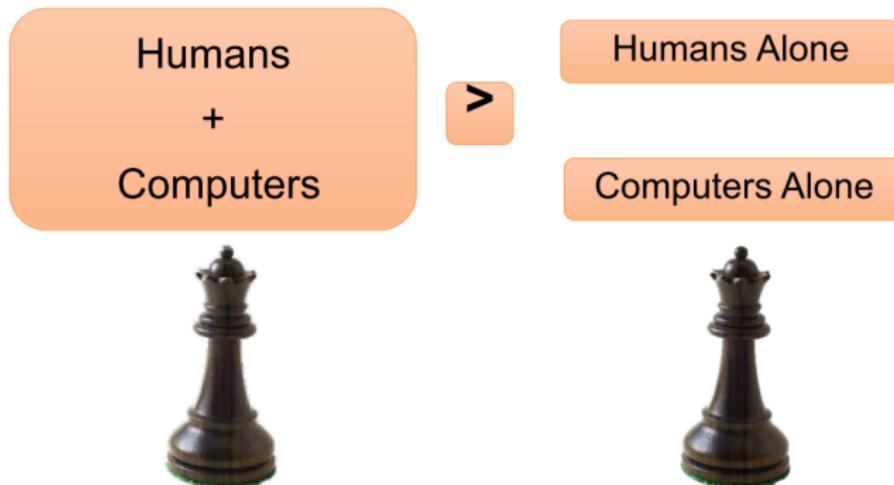
(Ref: What is Artificial Intelligence — Artificial Intelligence Tutorial For Beginners — Edureka)

## Fear: Are we being replaced?

- ▶ Yes. in tasks that are repetitive
- ▶ But not which require complex thinking and creativity

# Mostly

Technology Enhancing (Not Replacing) Humans



(Ref: "Artificial Intelligence Overview" - Harry Surden )

## Limits on Artificial Intelligence

- ▶ Many things still beyond the realm of AI
- ▶ No thinking computers
- ▶ No Abstract Reasoning
- ▶ Often AI systems Have Accuracy Limits
- ▶ Many things difficult to capture in data
- ▶ Sometimes Hard to interpret Systems

# Introduction to Generative AI

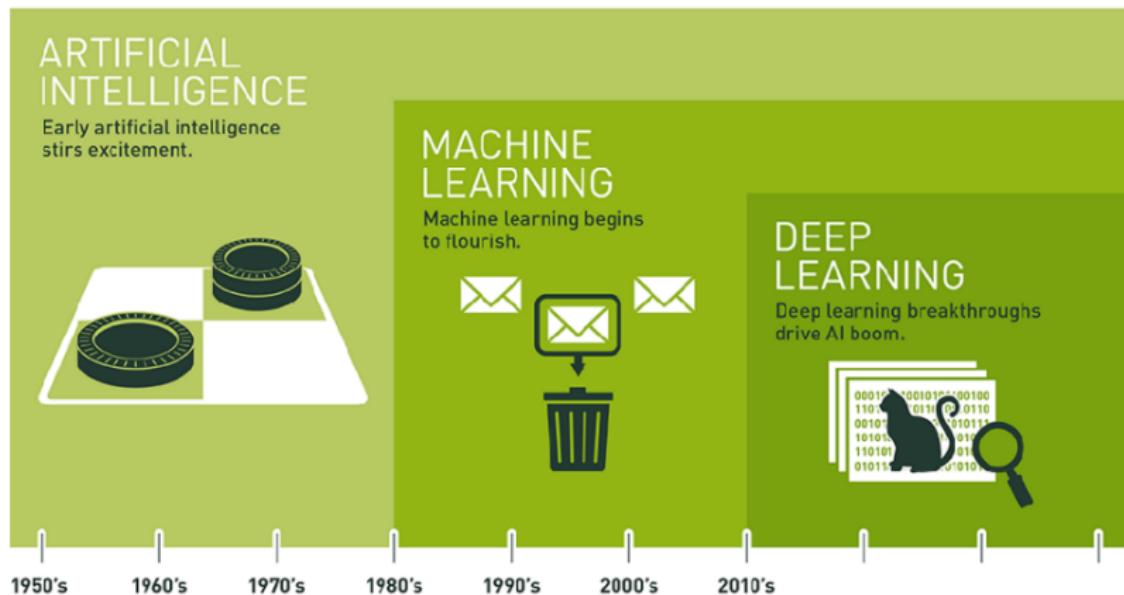
YHK

# Introduction

- ▶ What is Generative AI?
- ▶ What is not Generative AI?
- ▶ How is it related to AI-ML-DL?

# Relationship between AI, ML, DL

First, let's see what's AI-ML-DL and relationship among them.



(Ref: <https://blogs.nvidia.com/blog/2016/07/29/whats-difference-artificial-intelligence-machine-learning-deep-learning-ai/>)

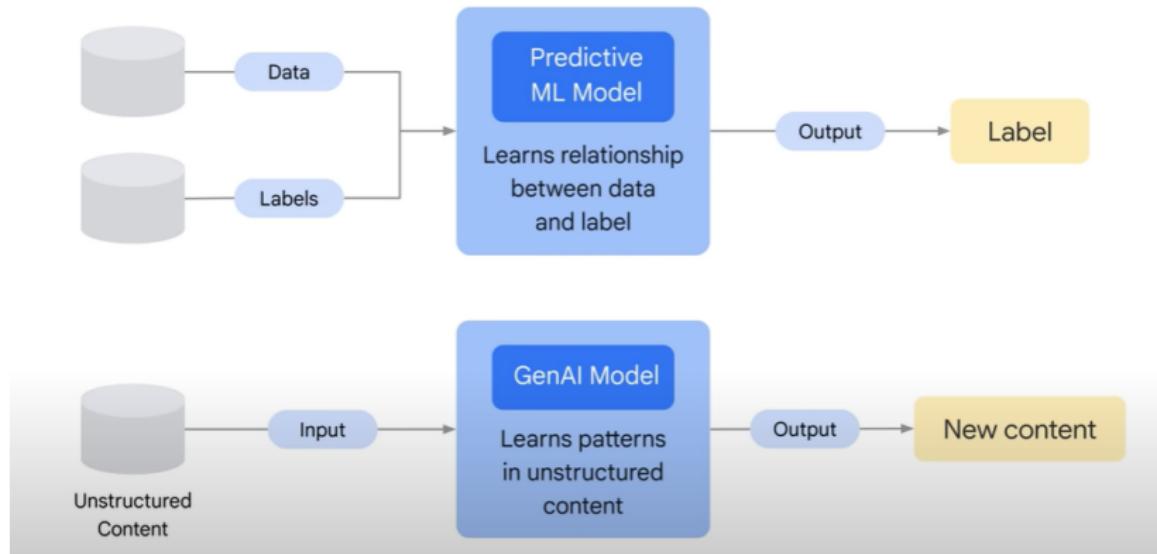
# What is Gen AI wrt AI, ML, DL

**Generative AI**  
is a **subset of**  
**Deep Learning**



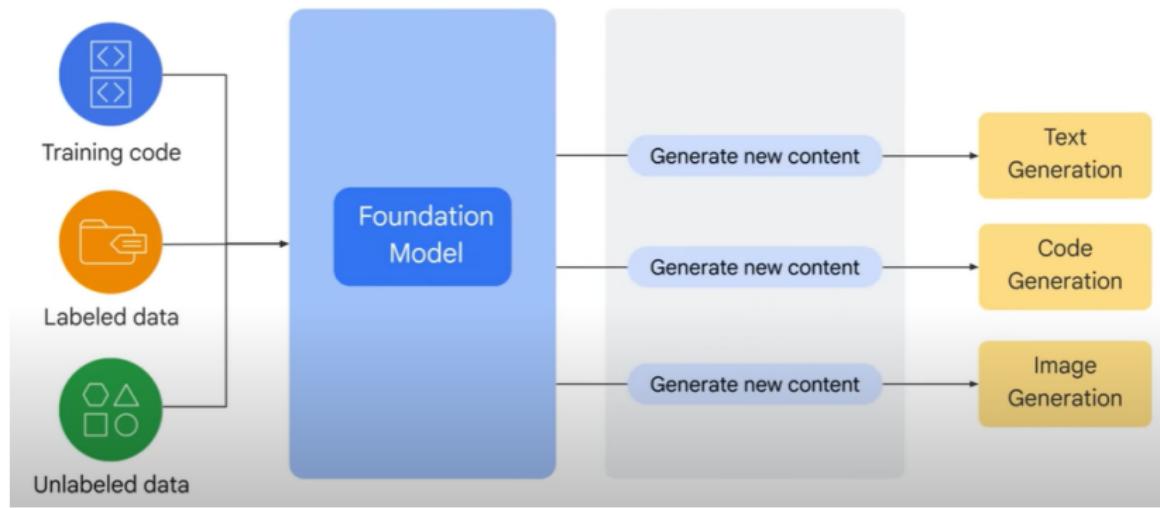
(Ref: Introduction to Generative AI - Google Cloud Tech)

## Types of Approaches



(Ref: Introduction to Generative AI - Google Cloud Tech)

# What is Foundation Model?



(Ref: Introduction to Generative AI - Google Cloud Tech)

## Same Problem, using different Technologies

YHK

## Difference across technologies, old to new

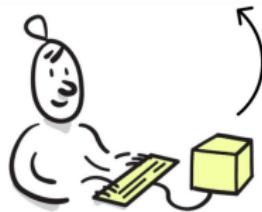
Lets see how the solutions to the problem of detecting a cat from images using traditional programming, deep learning, and generative AI, respectively.



## Traditional Programming

- ▶ Traditional programming involves writing explicit rules to detect a cat in images.
- ▶ Features like color, texture, and shape can be used to define these rules.
- ▶ However, designing accurate rules for complex patterns like cat detection can be challenging.
- ▶ It requires extensive domain knowledge and might not generalize well to different images.

```
cat:  
  type: animal  
  legs: 4  
  ears: 2  
  fur: yes  
  likes: yarn, catnip
```

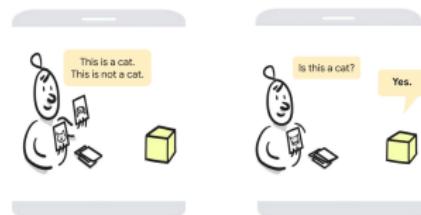


(Ref: Primer on LLM and Gen AI - Google Cloud)

YHK

# Deep Learning

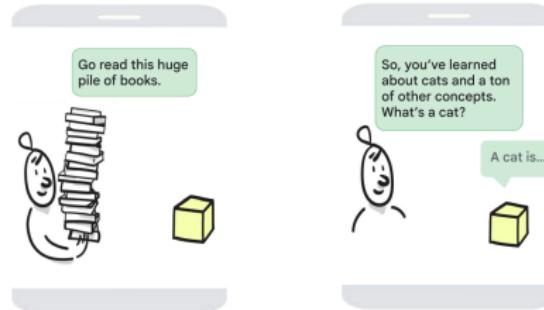
- ▶ Deep learning utilizes neural networks to automatically learn features for cat detection.
- ▶ Convolutional Neural Networks (CNNs) are particularly effective for image classification tasks.
- ▶ Large labeled datasets of cat images are used to train the network.
- ▶ The network learns to identify unique cat features and generalize them to detect cats in new images.
- ▶ Deep learning offers better accuracy and can handle complex patterns without explicit rule definition.



(Ref: Primer on LLM and Gen AI - Google Cloud)

# Generative AI

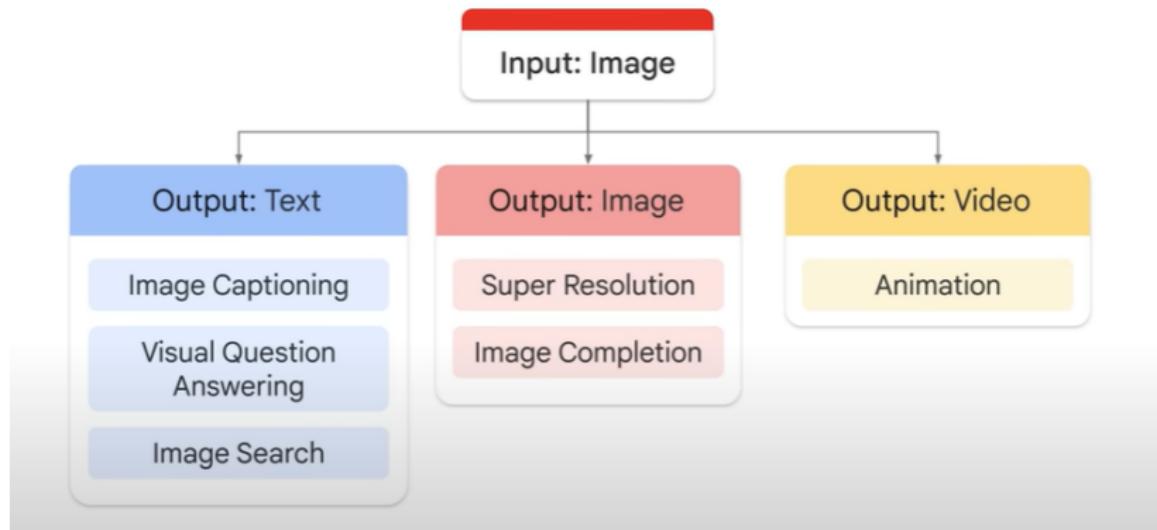
- ▶ Generative AI focuses on generating new data, including images of cats.
- ▶ Generative Adversarial Networks (GANs) are used to generate realistic cat images.
- ▶ The GAN consists of a generator and a discriminator that compete against each other.
- ▶ The generator learns to generate increasingly realistic cat images, while the discriminator learns to distinguish real from generated images.
- ▶ The generated cat images can be used to augment datasets for cat detection models.



(Ref: Primer on LLM and Gen AI - Google Cloud)

YHK

# Modalities in Generative AI

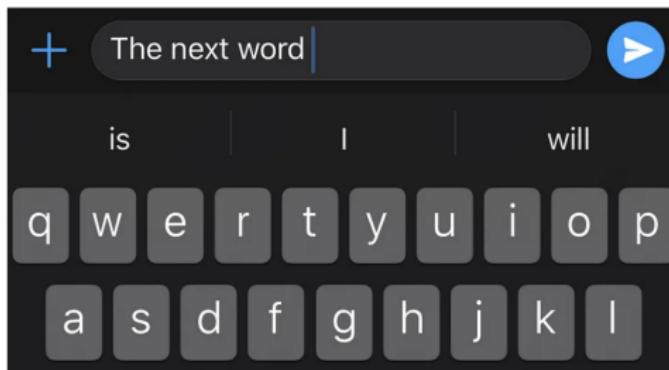


(Ref: Introduction to Generative AI - Google Cloud Tech)

Let's focus on the most popular modality ...

## What is a Language Models?

- ▶ While typing SMS, have you seen it suggests next word?
- ▶ While typing email, have you seen next few words are suggested?
- ▶ How does it suggest? (suggestions are not random, right?)
- ▶ In the past, for “Lets go for a . . .”, if you have typed ‘coffee’ 15 times, ‘movie’ say 4 times, then it learns that. Machine/Statistical Learning.
- ▶ Next time, when you type “Lets go for a ”, what will be suggested? why?
- ▶ This is called Language Model. Predicting the next word. When done continuously, one after other, it spits sentence, called Generative Model.



Next word prediction using language modeling in keyboards(Mandar Deshpande)

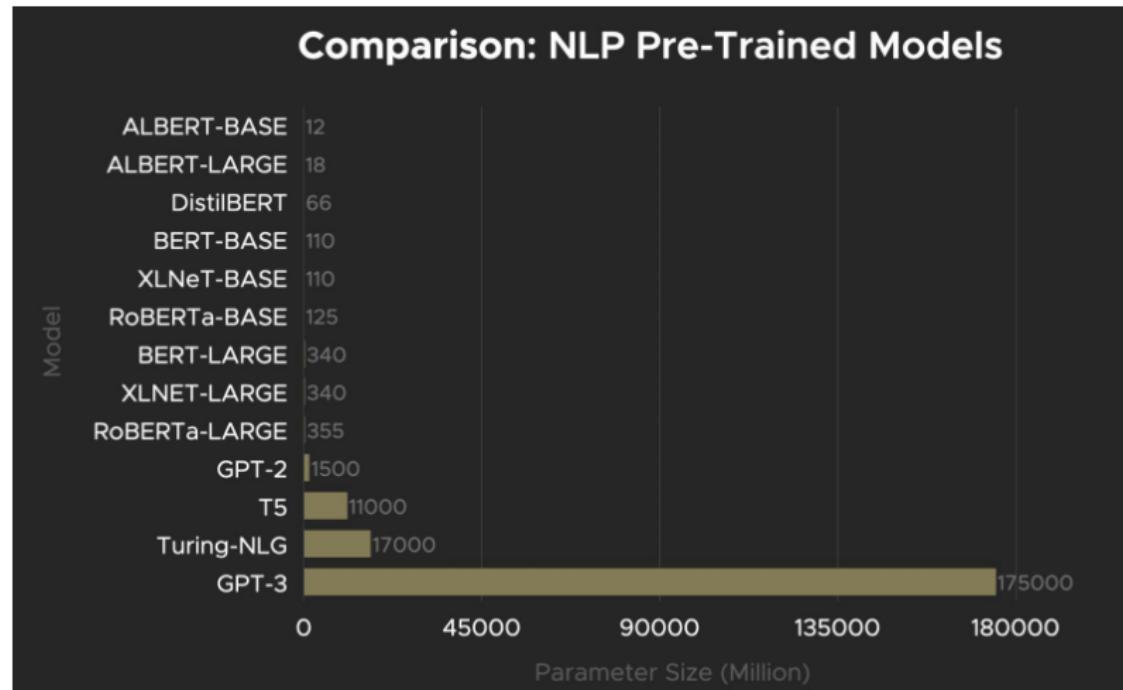
YHK

## Why they are called Large? Corpus

### GPT: Generative Pre-trained Transformers

- ▶ GPT-1 is pre-trained on the BooksCorpus dataset, containing 7000 books amounting to 5GB of data
- ▶ GPT-2 is pre-trained using the WebText dataset which is a more diverse set of internet data containing 8M documents for about 40 GB of data
- ▶ GPT-3 uses an expanded version of the WebText dataset, two internet-based books corpora that are not disclosed and the English-language Wikipedia which constituted 600 GB of data

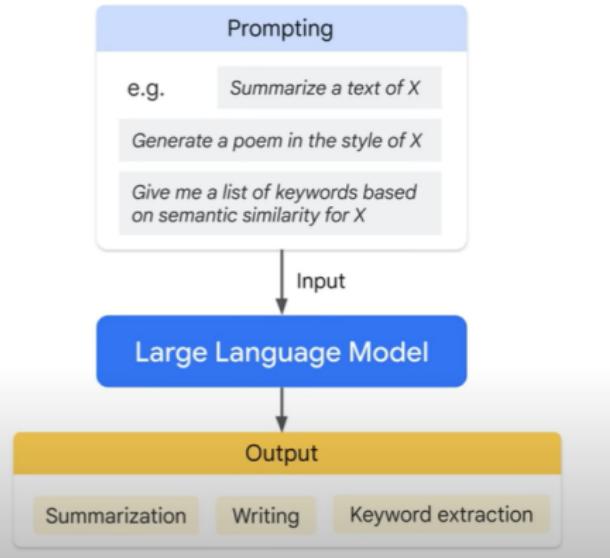
## Why they are called Large? Parameters



(Ref: Deus.ai <https://www.deus.ai/post/gpt-3-what-is-all-the-excitement-about>)

# Prompts driving Generative AI

**Prompt Design:**  
the quality of the  
input **determines** the  
quality of the output.



(Ref: Introduction to Generative AI - Google Cloud Tech)

# Making of a Sandwich



## Basic

Explain how to make a peanut butter and jelly sandwich



## Adding Roles

As a chef, explain to your assistant how to make a peanut butter and jelly sandwich



## Adding Constraints

Make a nut-free version of the sandwich due to a customer's nut allergy



## Adding Examples

Create two unique variations of the classic sandwich. Banana Nut Crunch: . . .

Triple Berry Blast: . . .



## Adding Contextual Information

As the head chef at 'The Sandwich Haven,' guide your new assistant to create specials for the menu



## Incorporating Feedback

Improve the sandwich based on customer feedback for less sweetness and a creative twist

## Time Constraints and Prioritization

Prepare an alternative fruit version for testing within a tight deadline



## Incorporating Multidisciplinary Knowledge

Use food presentation and garnishing techniques for a visually appealing sandwich



## Addressing Dietary Preferences

Prepare a vegan version using plant-based alternatives for all ingredients

## Reflection and Iteration

Reflect on feedback and iteratively refine the sandwich for better taste and appeal



## Self-Criticism

Explain how to make a peanut butter and jelly sandwich. Please re-read your above response. Any mistakes? If so, please identify and make the necessary edits.



## Chain-of-Thought

Explain how to make a peanut butter and jelly sandwich. Let's think step by step.

## Self-Consistency

Here are recipes of multiple sandwiches. Sandwich 1: recipe 1. Sandwich 2: recipe 2. .... Explain how to make a peanut butter jelly sandwich.

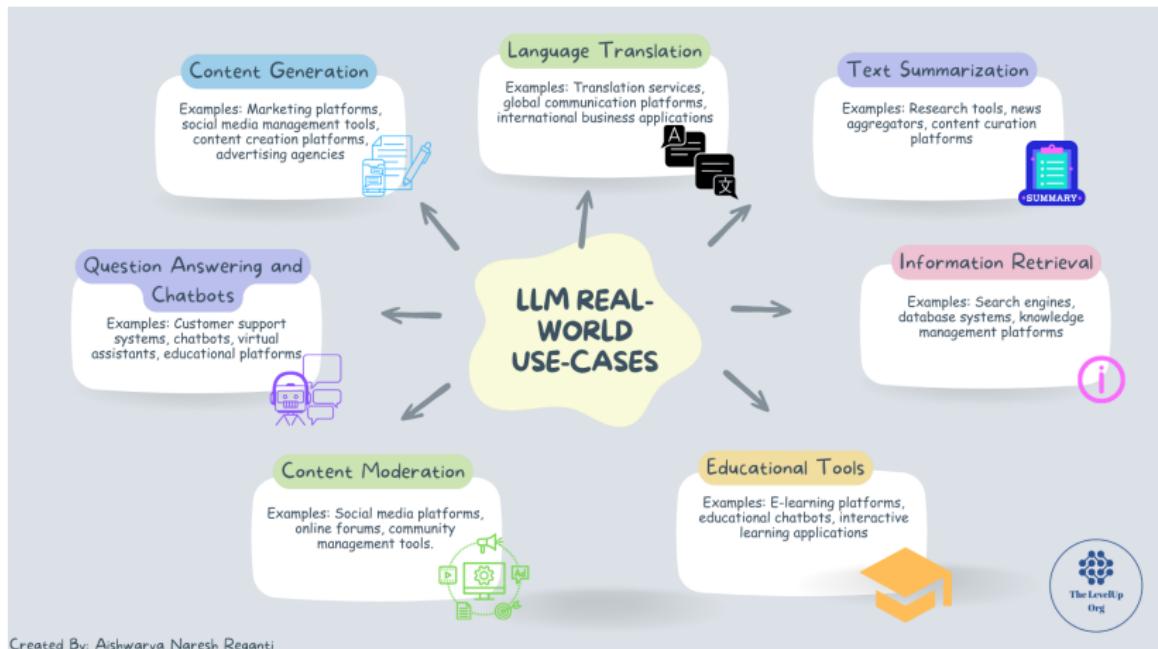


## Conclusions

(Ref: 2023 Kaggle AI Report on Generative AI, by Trushant Kalyanpur)



# LLM Real World Use Cases



Created By: Aishwarya Naresh Reganti

(Ref: Applied LLMs Mastery 2024 - Aishwarya Reganti)



# Advancements in Text-to-Image Generative AI

- ▶ Notable models: DALL-E/DALL-E 2, Midjourney, Stable Diffusion
- ▶ Creative expression, streamlined design
- ▶ Realistic, high-quality image generation
- ▶ Concerns: misuse, ethical implications
- ▶ Deepfakes, synthetic media for misinformation
- ▶ Risk of generating illegal, toxic content
- ▶ Challenges in ethical responsibility, moderation

# Copilots: Revolutionizing Coding

- ▶ AI assistance for software developers
- ▶ 92% programmers use AI tools (Github survey)
- ▶ Copilot users 55% faster in project completion
- ▶ Potential \$1.5 trillion GDP boost (productivity study)
- ▶ AI tools like Copilot enhance speed, efficiency
- ▶ Fewer errors in coding processes

# Industry Giants' Role in Text Generative AI

- ▶ Google, Meta, OpenAI - Pivotal contributions
- ▶ PaLM, Galactica, ChatGPT, GPT4 releases
- ▶ ChatGPT: Turning point in Text Generative AI
- ▶ LLMs for content creation, writing, storytelling
- ▶ Analyzing, organizing large textual data
- ▶ Efficient search engines, knowledge systems

# Advancements in LLM Accessibility

- ▶ Meta's LLaMa: Open-source alternatives to ChatGPT
- ▶ Google's Bard: AI chatbot response to ChatGPT
- ▶ QLoRA: Fine-tuning LLMs on consumer GPUs
- ▶ Broadening access to advanced technology
- ▶ LLMs empower diverse applications
- ▶ LLMs: Bridging the gap between innovation and accessibility



# Conclusions

YHK

# Progression

Models for prediction:

- ▶ On data, derive features, put statistical techniques like regression. One model per task. That's Machine Learning.
- ▶ Feed raw data, employ neural networks. One model per task. That's Deep Learning.
- ▶ Use Text data, get embeddings, use ML/DL, say for classification. One model per task. That's Natural Language Processing.
- ▶ Train neural network on large corpus, store weights and architecture, then add final layers for say classification on custom data+labels. That's Pretrained model. One model, many tasks.
- ▶ Train Large Language Model, just supply instructions on what to do, works. One model many tasks. Zero-shot, few-shots.

(More info at SaaS LLM <https://medium.com/google-developer-experts/saasgpt-84ba80265d0f>)

# New Programming Language?



Andrej Karpathy ✅  
@karpathy

...

The hottest new programming language is English

1:44 AM · Jan 25, 2023 · 1.9M Views

---

2,050 Retweets 284 Quote Tweets 17.9K Likes

---

(Ref: Prompt Engineering Sudalai Rajkumar)

# Summary of Prompt Engineering

## ► **Definition of Prompts**

- ▶ Prompts are initial text inputs provided to a model.
- ▶ Used by the model to generate responses or accomplish tasks.

## ► **Role of Prompts**

- ▶ Sets of instructions for AI or chatbots (e.g., ChatGPT).
- ▶ Applied in various tasks, including summarization, arithmetic problem-solving, and question-answering.

## ► **Objective of Prompt Engineering**

- ▶ Goal: Refine prompts to enhance model accuracy and relevance in outputs.
- ▶ Central to improving the performance of language models.

## ► **Prevalent Prompt Types**

- ▶ Various prompt types exist, with a focus on two widely used methodologies:
  - ▶ Zero-shot prompting
  - ▶ Few-shot prompting

# ChatGPT Ultimate Prompting Guide

- ▶ Tone: Specify the desired tone (e.g., formal, casual, informative, persuasive).
- ▶ Format: Define the format or structure (e.g., essay, bullet points, outline, dialogue).
- ▶ Act as: Indicate a role or perspective to adopt (e.g., expert, critic, enthusiast).
- ▶ Objective: State the goal or purpose of the response (e.g., inform, persuade, entertain).
- ▶ Context: Provide background information, data, or context for accurate content generation.
- ▶ Scope: Define the scope or range of the topic.
- ▶ Keywords: List important keywords or phrases to be included.
- ▶ Limitations: Specify constraints, such as word or character count.
- ▶ Examples: Provide examples of desired style, structure, or content.
- ▶ Deadline: Mention deadlines or time frames for time-sensitive responses.

(Ref: LinkedIn post by Generative AI, Twitter by Aadit Sheth, Source : Reddit)



# ChatGPT Ultimate Prompting Guide

- ▶ Audience: Specify the target audience for tailored content.
- ▶ Language: Indicate the language for the response, if different from the prompt.
- ▶ Citations: Request inclusion of citations or sources to support information.
- ▶ Points of view: Ask the AI to consider multiple perspectives or opinions.
- ▶ Counter arguments: Request addressing potential counterarguments.
- ▶ Terminology: Specify industry-specific or technical terms to use or avoid.
- ▶ Analogies: Ask the AI to use analogies or examples to clarify concepts.
- ▶ Quotes: Request inclusion of relevant quotes or statements from experts.
- ▶ Statistics: Encourage the use of statistics or data to support claims.
- ▶ Visual elements: Inquire about including charts, graphs, or images.
- ▶ Call to action: Request a clear call to action or next steps.
- ▶ Sensitivity: Mention sensitive topics or issues to be handled with care or avoided.

(Ref: LinkedIn post by Generative AI, Twitter by Aadit Sheth, Source : Reddit)



## Interaction Guidelines: Avoid Misuses

- ▶ Factual Accuracy: Interactions must be free from factual inaccuracies that can be challenged by social media or journalists.
- ▶ Negative Debates: Avoid discussing topics that fuel negative or concerning online debates, such as AI sentience, AI in education, AI-driven job displacements, and politically divisive issues.
- ▶ Minors' Involvement: Do not include use cases specifically targeting or involving individuals under 18 years old.
- ▶ Sensitivity and Misinformation: Prevent the inclusion of sensitive, misleading, or hazardous responses.
- ▶ Search and Google Assistant: Interactions that require basic, straightforward answers are better suited for Search or Google Assistant.
- ▶ Financial/Legal/Medical Advice: Refrain from providing advice related to financial matters, legal issues, or medical concerns.
- ▶ Brand Names and Trademarks: Avoid mentioning specific brand names, trademarks, or public figures (except historical figures).
- ▶ No Reviews or Tweets: Do not request reviews of restaurants, businesses, or tweets to minimize the risk of associating with bots.
- ▶ Avoid Personification: Refrain from personifying the product or brand and from encouraging users to address Bard by name.

## Limitations

Boie is a real company, the product name is not real. So, see what you get ...

```
1 prompt = f"""
2 Tell me about AeroGlide UltraSlim Smart Toothbrush by Boie
3 """
4
5 response = get_completion(prompt)
6 print(response)
```

# What Next?

YHK

# The Career of the Future

## Software 3.0

---

“Programming [is] moving from curating datasets to curating prompts to make the meta learner “get” the task it’s supposed to be doing.”

Source: [@karpathy](#)



(Ref: The Complete Prompt Engineering for AI Bootcamp (2023))

## New Roles?

Coming up with good prompt is a combination of art and science



Alexandr Wang @alexandr\_wang

...

Today, [@goodside](#) joined [@scale\\_AI](#) as a Staff Prompt Engineer.

I am going to assert that Riley is the first Staff Prompt Engineer hired \*anywhere\*.

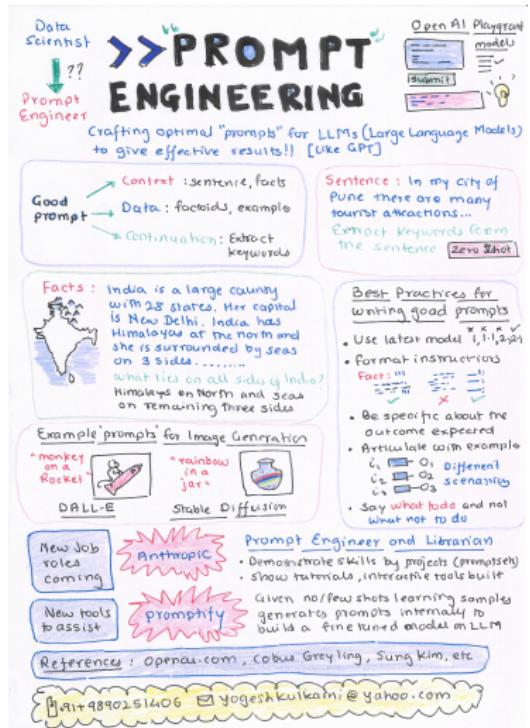
(Ref: Prompt Engineering Sudalai Rajkumar)

## Read on to learn how to engineer good prompts!

- ▶ Shin, T., Razeghi, Y., Logan IV, R. L., Wallace, E., & Singh, S. (2020). AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP).  
<https://doi.org/10.18653/v1/2020.emnlp-main.346>
- ▶ Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., & Iwasawa, Y. (2022). Large Language Models are Zero-Shot Reasoners.
- ▶ Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., & Neubig, G. (2022). Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. ACM Computing Surveys.  
<https://doi.org/10.1145/3560815>
- ▶ Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., ... Amodei, D. (2020). Language Models are Few-Shot Learners.
- ▶ Zhao, T. Z., Wallace, E., Feng, S., Klein, D., & Singh, S. (2021). Calibrate Before Use: Improving Few-Shot Performance of Language Models.



# My Sketchnote



(Ref: <https://medium.com/technology-hits/prompting-is-all-you-need-5dddb82bd022>)

## Take Aways

Prompt Engineering is an Iterative Process:

- ▶ Try something
- ▶ Analyze where the results do not match the expectations
- ▶ Clarify instructions, gives examples, specify output format, specify constraints, etc
- ▶ Test on a batch of known results.

## Quality of Ideas

"I don't think we'll still be doing prompt engineering in five years [i.e.] figuring out how to hack the prompt by adding one **magic word** to the end that changes everything else."

"What will always matter is the **quality of ideas** and the understanding of what you want."

Source: [Sam Altman](#)

(Ref: The Complete Prompt Engineering for AI Bootcamp (2023))



## Resources

- ▶ Prompt Engineering Guide  
<https://github.com/dair-ai/Prompt-Engineering-Guide>
- ▶ Awesome ChatGPT Prompts  
<https://github.com/f/awesome-chatgpt-prompts/>
- ▶ ChatGPT Prompt Engineering for Developers - Deep Learning AI
- ▶ Learn Prompting <https://learnprompting.org/docs/intro>
- ▶ Types of Prompts with Practical examples - Dr. Naveed Siddiqui
- ▶ AI Prompt Database  
<https://justunderstandingdata.notion.site/d98dcc9a6736471584d53cc8b2a5c30d?v>

## References

- ▶ Introduction to Generative AI - Google Cloud Tech
- ▶ Generative AI Presentation - Laura Worden

## Newsletters to subscribe

- ▶ **The Batch by DeepLearning.AI:**
  - ▶ Summarizes diverse AI news with nuanced viewpoints.
  - ▶ Andrew Ng's thought leadership adds significant value.
- ▶ **The Rundown AI by Rowan Cheung:**
  - ▶ Go-to for generative AI events and product innovations.
  - ▶ Quick rundown with bullet point details for easy comprehension.
- ▶ **AI Supremacy by Michael Spencer:**
  - ▶ Personal writing style with in-depth exploration.
  - ▶ Offers multiple perspectives on AI topics.
- ▶ **Ahead of AI by Sebastian Raschka, PhD:**
  - ▶ Technical focus covering applied deep learning and generative AI.
  - ▶ Valuable insights for those seeking in-depth technical content.
- ▶ **To Data and Beyond by Youssef Hosni:**
  - ▶ Resource hub for hands-on projects, learning roadmaps, and research papers.
  - ▶ Ideal for those looking to dive into practical aspects of AI.

Thanks ...

- ▶ Search "**Yogesh Haribhau Kulkarni**" on Google and follow me on LinkedIn and Medium
- ▶ Office Hours: Saturdays, 2 to 3 pm (IST); Free-Open to all; email for appointment.
- ▶ Email: yogeshkulkarni at yahoo dot com



(<https://www.linkedin.com/in/yogeshkulkarni/>)



(<https://www.linkedin.com/in/yogeshkulkarni/>)



(<https://www.github.com/yogeshhk/>)

YHK