

Introduction to Large Language Models (LLM)

*By Raghava
3rd Year,
M.Sc Artificial Intelligence and Machine Learning*

Language Model (LM)

A language model is a computer program that is trained on a large amount of text data to generate language outputs that are coherent and natural-sounding. It predicts the next word or character in a sequence of text, based on the context and patterns learned from the training data.

Input: "The cat sat on the _____"

Language Model Output: "mat/floor/etc..."

Explanation: The language model predicts the next word "mat" based on the context of the sentence and the patterns learned from the training data.

Some popular language models include BERT and Transformers.

What is LLM?

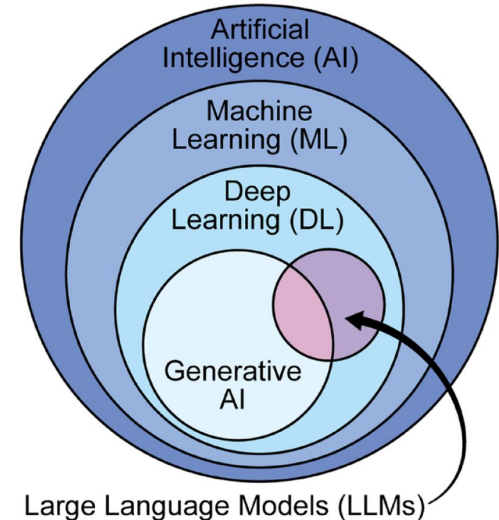
LLM stands for Large Language Model. A Large Language Model is a type of artificial intelligence (AI) model that is trained on a vast amount of text data to generate language outputs that are coherent and natural-sounding. LLMs are designed to process and understand human language, and they have numerous applications, possibly every domain.

Popular LLMs:

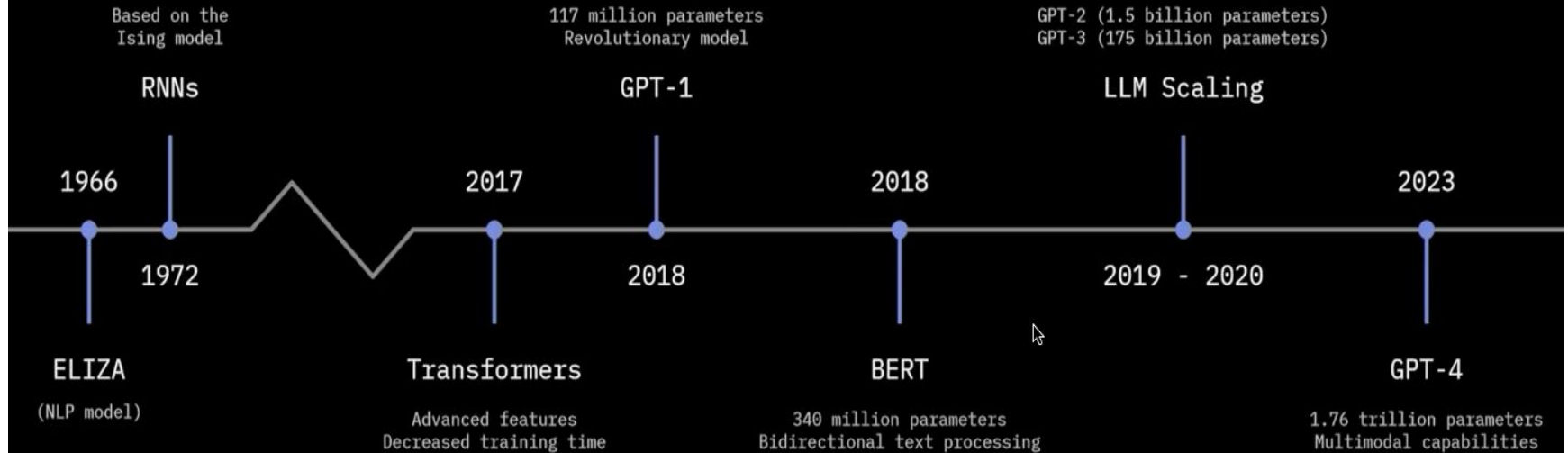
ChatGPT by OpenAI

Gemini & Gemma by Google

llama by Meta



History of Large Language Models



Foundational models

Foundational models in large language models (LLMs) refer to the initial models that serve as the basis or foundation for further research, development, and fine-tuning within the field of natural language processing (NLP). These models are typically pre-trained on vast amounts of text data using unsupervised learning techniques, such as self-supervised learning or autoencoding, to learn the underlying patterns and structures of language.

GPT (Generative Pre-trained Transformer) developed by OpenAI

BERT (Bidirectional Encoder Representations from Transformers) developed by Google

T5 (Text-To-Text Transfer Transformer) developed by Google Brain

LLM Modalities

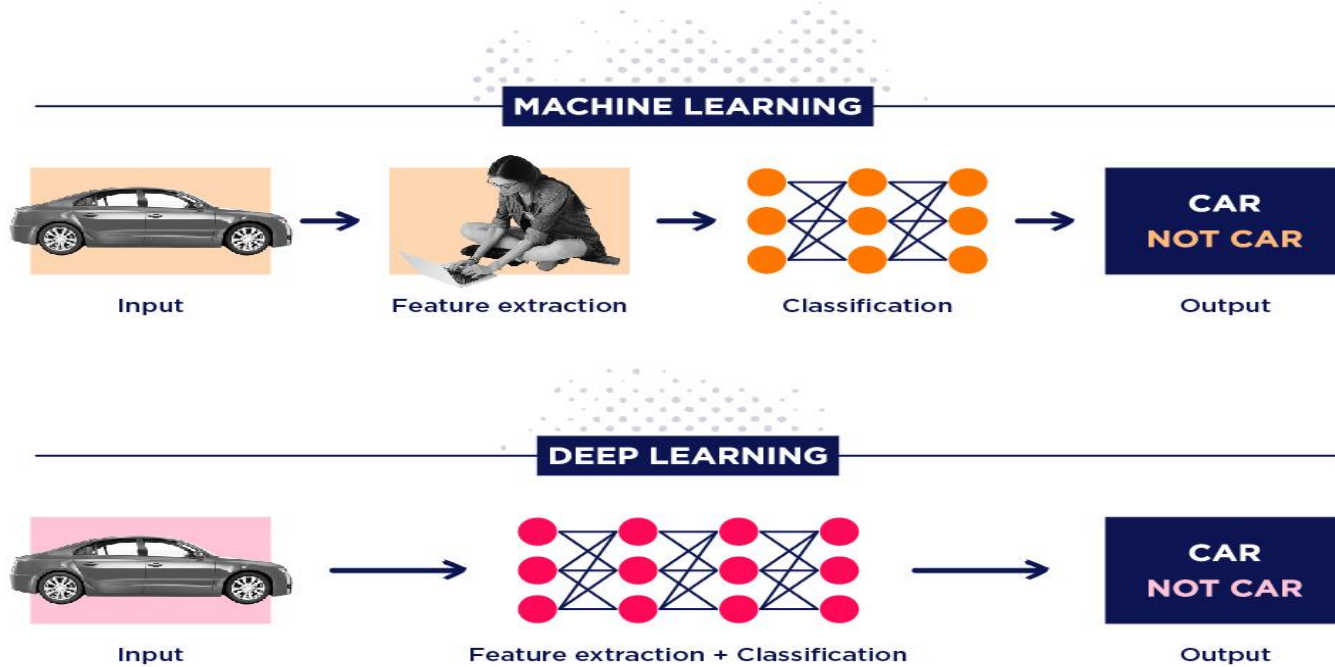
TEXT MODELS: These models focus primarily on processing and generating text data. They are trained on large text corpora and excel at tasks such as language modeling, text generation. *GPT (Generative Pre-trained Transformer) models developed by OpenAI and BERT (Bidirectional Encoder Representations from Transformers) developed by Google.*

IMAGE-TEXT MODELS: LLMs that integrate both text and image modalities can understand and generate content based on both textual and visual inputs. *Gemini developed by Google.*

TEXT-VIDEO MODELS: Given a text prompt as a input, these models can generate video based on the context from the prompt and make a new real life video. *Sora by OpenAI*

MULTIMODAL MODELS: These models combine information from multiple modalities (text, image, audio) to generate coherent outputs that leverage the strengths of each modality. *Google AI's PaLM-E, OpenAI's GPT-4V*

Training ML vs DL



Types of LLMs

AUTOREGRESSIVE LLMS: Generate text one token at a time, predicting the next token based on previous context. Examples: GPT-3, GPT-4.

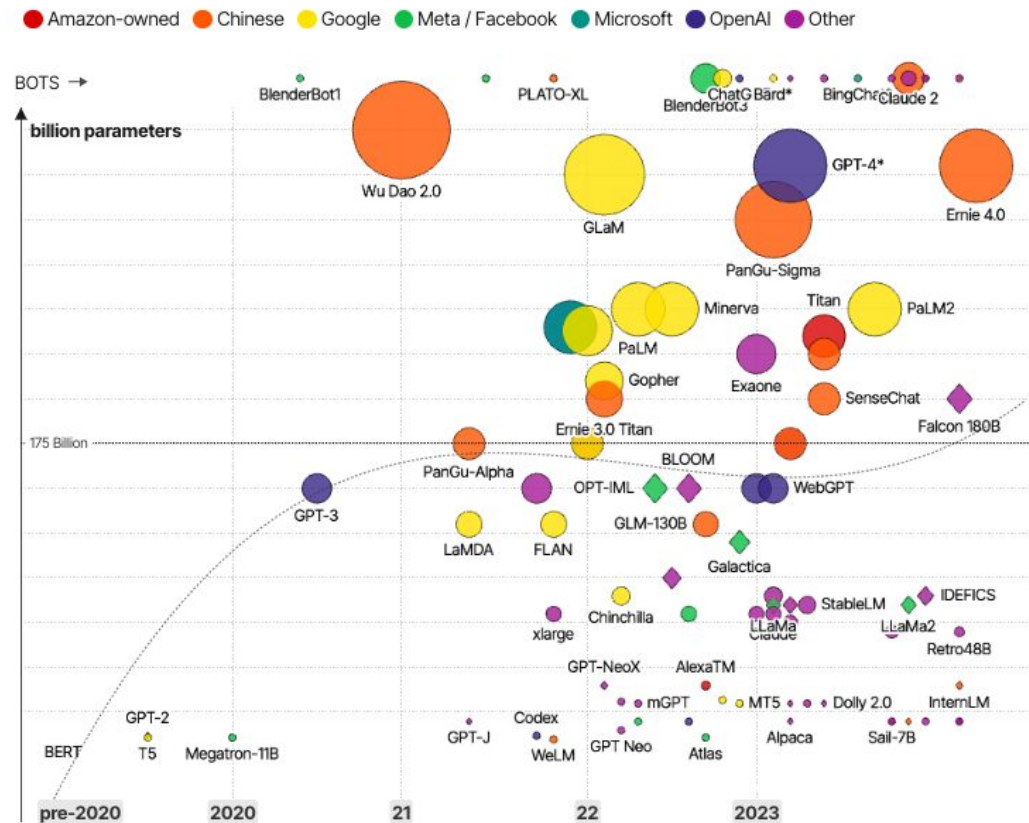
AUTOENCODER LLMS: Compress input text into a smaller representation and then reconstruct the original text. Examples: BERT, RoBERTa.

DECODER-ONLY LLMS: Use only the decoder part of the transformer architecture, generating text without encoding input text. Examples: GPT-3, GPT-4.

ENCODER-DECODER LLMS: Use both encoder and decoder parts of the transformer architecture, encoding input text and generating output text. Examples: BERT, RoBERTa.

More...

The Rise and Rise of A.I. Large Language Models (LLMs) & their associated bots like ChatGPT



Source: Datacamp

David McCandless, Tom Evans, Paul Barton
Information is Beautiful // UPDATED 2nd Nov 2023

source: news reports, [LifeArchitect.ai](https://lifeaiarchitect.ai)
* = parameters undisclosed // see [the data](https://the-data)

Buzz words in the world of LLMs

PARAMETERS:

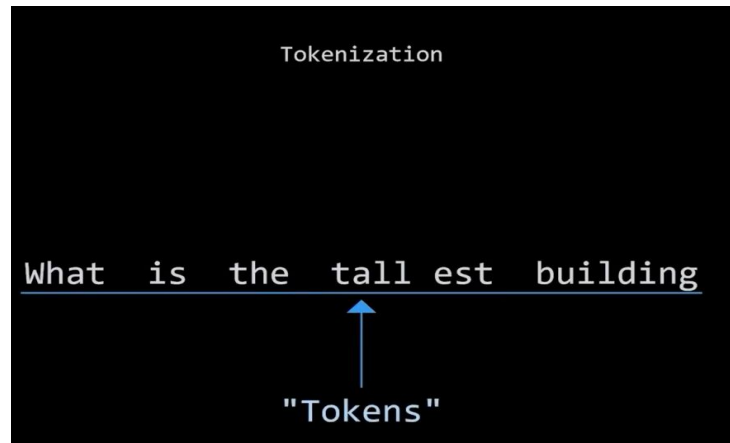
The purpose of having a large number of parameters in an LLM is to increase its capacity to learn complex patterns and relationships in language. More parameters allow the model to capture finer nuances of language, including syntax, semantics, context, and domain-specific knowledge. Larger models with more parameters also require more computational resources for training and inference.

Example: Gemma 2B, 7B

- A 7B (7 billion) parameter model has 7,000,000,000 learnable weights and biases.
- A 2B (2 billion) parameter model has 2,000,000,000 learnable weights and biases.

TOKENIZATION:

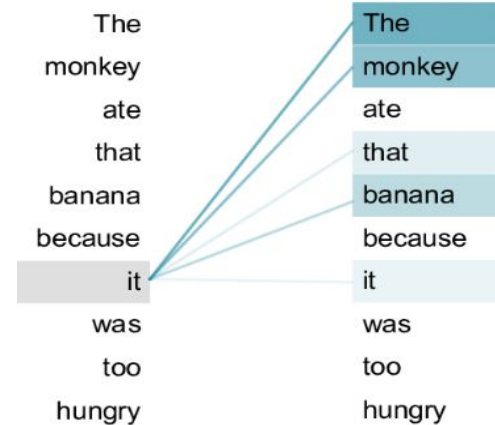
It is the process of breaking down a sequence of text into smaller units called tokens. These tokens can be words, subwords, characters, or any other meaningful units, depending on the specific tokenization strategy used. Tokenization is a fundamental step in natural language processing (NLP) tasks, as it allows machines to understand and process human language.



Buzz words in the world of LLMs

SELF ATTENTION:

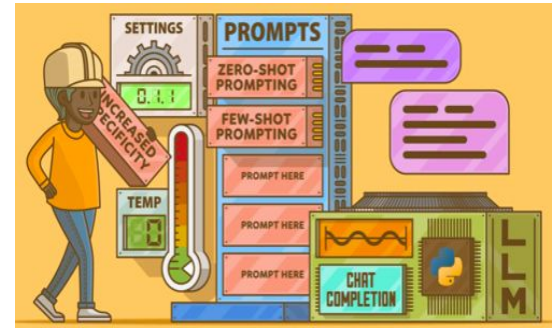
Self-attention is a mechanism that computes the importance of each token in a sequence by considering its relationships with all other tokens in the same sequence. It allows the model to assign different levels of importance or attention to each token based on its relevance to other tokens in the sequence.



PROMPT ENGINEERING:

Prompt engineering refers to the process of designing or crafting natural language prompts or instructions to control the behavior of large language models (LLMs) and elicit desired responses. It involves carefully constructing input text that guides the LLM to produce outputs aligned with specific objectives or tasks. *It includes Instructions, Examples, Constraints, Contextual Information.*

Example: Give me 300 words assignment content on methods of classification in machine learning with examples include KNN, K-Means etc.



Buzz words in the world of LLMs

HALLUCINATION

It refers to a phenomenon where the model generates text or outputs that are not grounded in the input data or are inconsistent with reality. Hallucinations occur when the model generates content that is either unrelated to the input prompt or contains factual inaccuracies, contradictions, or nonsensical information.

Confidently incorrect: The model produces answers that are wrong, but with high confidence, as if it believes they are correct.

Not supported by training data: The responses are not based on any actual training data or evidence, but rather the model's own internal workings.

Often plausible-sounding: Hallucinated responses can be grammatically correct, coherent, and even persuasive, making them difficult to detect.

Example:
What is crane?



LLM Leaderboard

Tracks open-source LLMs and chatbots: It focuses on models freely available for everyone to use and improve.

Ranks models based on benchmarks: It evaluates models on various tasks like question answering and summarization.

Provides detailed results: You can explore how each model performs on individual benchmarks.

Collection of top models: It features a curated list of high-performing LLMs.

The screenshot shows the 'Open LLM Leaderboard' interface on HuggingFace Spaces. At the top, it indicates the space is by HuggingFaceH4, has 9.39k likes, and is running on CPU with an upgrade button. The main title is 'Open LLM Leaderboard'. Below this, there are navigation links: 'LLM Benchmark', 'Metrics through time', 'About', 'FAQ', and 'Submit'. A search bar allows users to search for models or licenses. There are several filter sections: 'Select columns to show' with checkboxes for Average, ARC, HellaSwag, MMLU, TruthfulQA, Winogrande, GSM8K, Type, Architecture, Precision, Merged, Hub License, #Params (B), Hub, and Model sha; 'Hide models' with checkboxes for Private or deleted, Contains a merge/merge, Flagged, and MoE; 'Model types' with checkboxes for pretrained, continuously pretrained, fine-tuned on domain-specific datasets, chat models (RLHF, DPO, IFT, ...), base merges and moerges, and a question mark; 'Precision' with checkboxes for float16, bfloat16, 8bit, 4bit, GPTQ, and a question mark; and 'Model sizes (in billions of parameters)' with checkboxes for a question mark, ~1.5, ~3, ~7, ~13, ~35, ~60, and 70+. The main table displays a list of models with their scores across various benchmarks. The table has columns for Model, Average, ARC, HellaSwag, MMLU, and Truthful.

Model	Average	ARC	HellaSwag	MMLU	Truthful
daavidkim285/Rhea-72b-v0.5	81.22	79.78	91.15	77.95	74.5
MTSAIR/MultiVerse_70B	81	78.67	89.77	78.22	75.18
MTSAIR/MultiVerse_70B	80.98	78.58	89.74	78.27	75.09
SF-Foundation/Ein-72B-v0.11	80.81	76.79	89.02	77.2	79.02
SF-Foundation/Ein-72B-v0.13	80.79	76.19	89.44	77.07	77.82
SF-Foundation/Ein-72B-v0.12	80.72	76.19	89.46	77.17	77.78
sharitsai/Smile-72B-v0.1	80.68	76.07	89.77	77.15	76.67

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

Evaluation Metrics for LLMs

ARC (question answering): Evaluates a model's ability to answer questions based on a given text.

HellaSwag (natural language understanding): Assesses a model's ability to understand natural language by identifying the correct ending to a story.

TruthfulQA (factual accuracy): Evaluates a model's ability to provide accurate and truthful answers to questions.

MMLU (massive multitask learning): Tests a model's ability to perform well across a wide range of tasks simultaneously.

ROUGE scores (summarization quality): Measures the quality of a model's summarization capabilities by comparing it to human-written summaries.

BLEU scores (machine translation quality): Evaluates the quality of a model's machine translation outputs by comparing them to human-translated references.

Applications of LLM

Chatbots and Virtual Assistants: Conversational AI, enabling chatbots to understand and respond to user queries.

Text Summarization: Summarize long texts, extracting key points and saving time.

Content Generation: Create content, such as articles, stories, and even entire books.

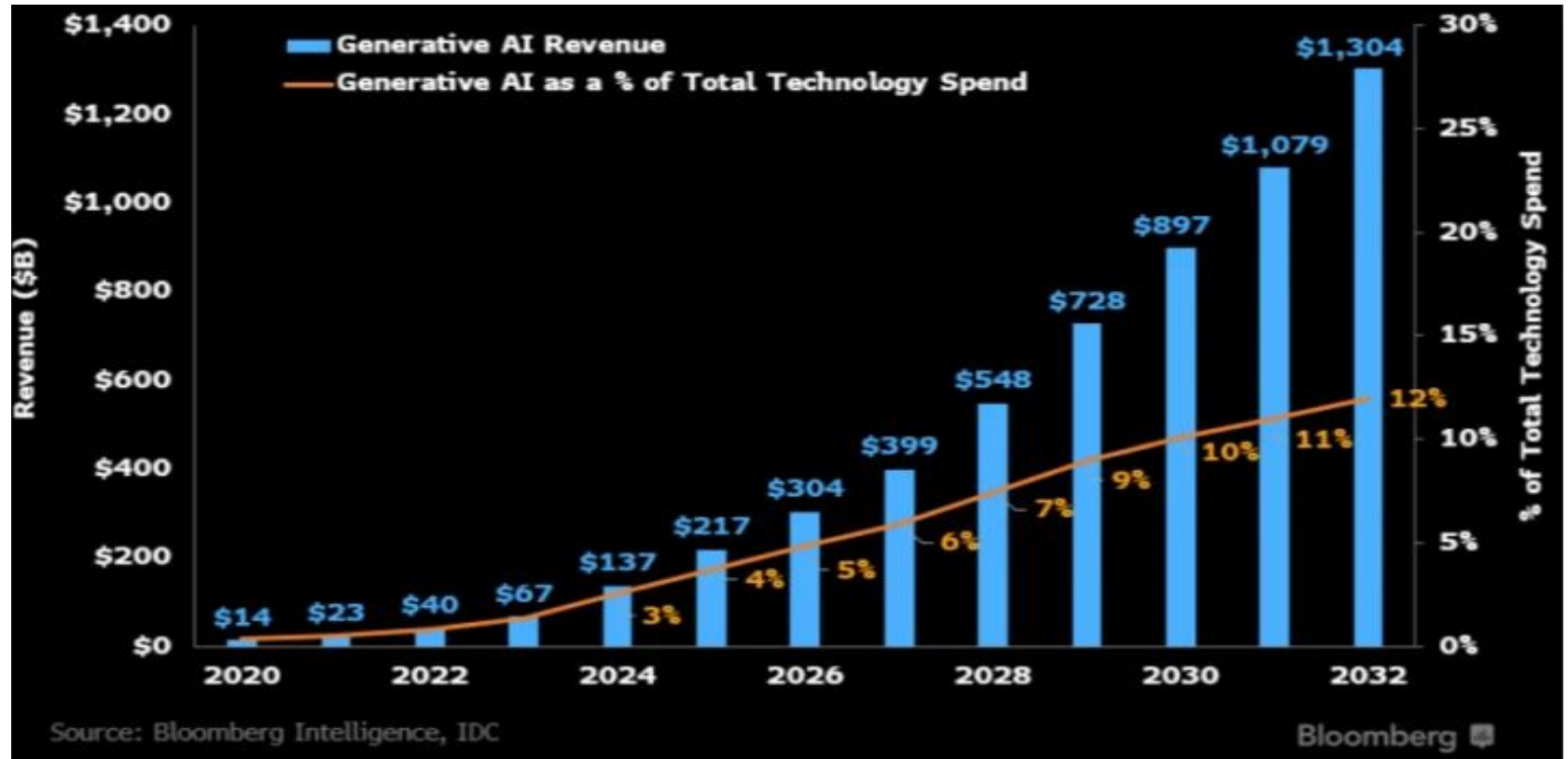
Question Answering: Answer questions, providing accurate and relevant information based on context.

Creative Writing: Assist writers with ideas, suggestions, and even co-authoring content.

Customer Service: Power automated customer support, providing quick and accurate solutions.

WHAT NOT ?! Much more ...

Gen AI in business



New Job Roles

Prompt Engineer

AI Trainer

AI Ethics Officer

Data Annotation Specialist

NLP Engineer

AI Researcher

Urban Planner

AI Personality Designer

These a just a few...



Quiz

Since the last quarter of year 2022, (ie. The launch & rise of LLMs, Generative AI such as ChatGPT, Gemini etc...).

Which company had a significant rise in their market share and net profit?

NVIDIA

To train and run models effectively. Nvidia's GPUs (Graphics Processing Units) are particularly well-suited for this task, due to their parallel processing architecture.

Nvidia A100 GPUs. Nvidia is ramping up its next-generation of GPUs, called H100, which delivers six times more throughput than the A100.



Resources to learn...

Youtube Channels:

Krish Naik: <https://www.youtube.com/@krishnaik06> (*For beginners*)

Deep Learning AI by Prof. Andrew NG : <https://www.youtube.com/@Deeplearningai>

CodeBasics: <https://www.youtube.com/@codebasics>

Assembly AI: <https://www.youtube.com/@AssemblyAI>

Varun Mayya: <https://www.youtube.com/@VarunMayya> (*Know new inventions*)

Nicholas Renotte: <https://www.youtube.com/@NicholasRenotte>

Resources to learn...

Websites:

Hugging Face: <https://huggingface.co/learn/nlp-course/chapter0/1>

Cohere: <https://docs.cohere.com/docs/llmu>

Datacamp: <https://www.datacamp.com/courses/large-language-models-llms-concepts>

Deep Learning AI: <https://www.deeplearning.ai/courses/>

Full Stack DL: <https://fullstackdeeplearning.com/llm-bootcamp/>

KDNuggets: <https://www.kdnuggets.com/a-comprehensive-list-of-resources-to-master-large-language-models>

Resources (BONUS)

You Can get free access to learn LLM & Generative AI in Google.

Google offers FREE BADGE Course for beginners and students

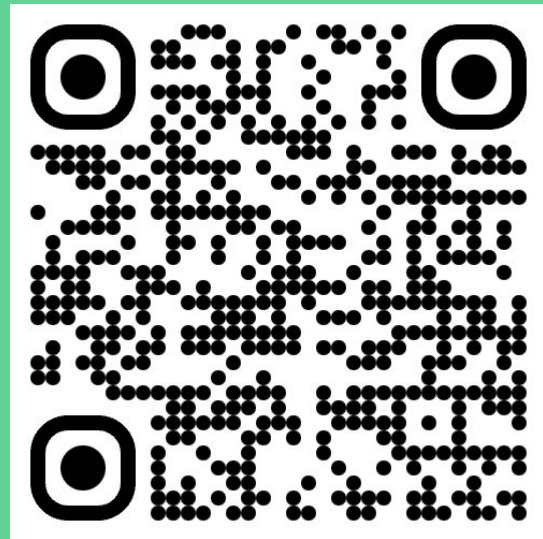
Checkout:

<https://www.cloudskillsboost.google/catalog>

The screenshot displays the Google Cloud Skills Boost interface. At the top, there's a navigation bar with links to Dashboard, Paths, Explore, Profile, and Subscriptions. A user's score of 0 pts is shown. Below the navigation bar, the page title is "Google Cloud Skills Boost". The main content area features three course cards, each with a "COMPLETION BADGE" icon.

- 01 Introduction to Generative AI**
 - Course icon: A blue checkmark inside a circle.
 - Duration: 45 minutes
 - Level: Introductory
 - Points: +800 pts
 - Description: This is an introductory level microlearning course aimed at explaining what Generative AI is, how it is used, and how it differs from traditional machine learning methods. It also covers Google Tools to help you develop your own Gen AI...
 - Status: ✓ Course complete
- 02 Introduction to Large Language Models**
 - Course icon: A blue checkmark inside a circle.
 - Duration: 30 minutes
 - Level: Introductory
 - Points: +800 pts
 - Description: This is an introductory level micro-learning course that explores what large language models (LLM) are, the use cases where they can be utilized, and how you can use prompt tuning to enhance LLM performance. It also covers Google tools to...
 - Status: ✓ Course complete
- 03 Introduction to Responsible AI**
 - Course icon: A blue checkmark inside a circle.
 - Duration: 30 minutes
 - Level: Introductory
 - Points: +800 pts
 - Description: This is an introductory-level microlearning course aimed at explaining what responsible AI is, why it's important, and how Google implements responsible AI in their products. It also introduces Google's 7 AI principles.
 - Status: → Start course

Thankyou



Connect with ME!
@raghavtwenty