

INT 426

INTRODUCTION TO GENAI

GENAI

- ↳ type of AI technology that can produce various types of content, including text, imagery, audio and synthetic data.

AI

- ↳ branch of CS that deals with the creation of intelligence agents which are systems that can reason, and learn and act autonomously.
- ↳ has to do with the theory & methods to build machines that think & act like humans

ML

- ↳ subfield of AI
- ↳ program or system that trains a model from input data.
- ↳ gives computer the ability to learn without explicit programming.
- ↳ Two of the most common classes of ML models
 - (1) Supervised models (labelled data)
 - (2) Unsupervised models (unlabelled data)
- ↳ LABELLED DATA data that comes with a tag like a name, a type or a number.
- ↳ UNLABELLED DATA data that doesn't come with a tag
- ↳ SUPERVISED LEARNING, the model learns from past examples to predict future values, ~~in this case~~
- ↳ UNSUPERVISED LEARNING, is all about discovery, about looking at the raw data and seeing if it naturally falls into groups.
- ↳ In supervised learning, testing values or x are input into the model. The model outputs a prediction & compares that prediction to the training data used to train the model. If the predicted test data values & actual training data values are far apart, that's called error. And the model tries to reduce this error until the predicted & actual are closer together. This is called CLASSIC OPTIMIZATION PROBLEM

DEEP LEARNING

- ↳ Subset of ML
- ↳ type of ML that uses artificial neural networks, allowing them to process more complex patterns than ML
- ↳ ARTIFICIAL NEURAL NETWORK are inspired by humans. They are made up of many inter connected nodes or neurons that can learn to perform tasks by processing data & making predictions.
- ↳ Neural networks can use both labelled & unlabelled data. This is called SEMI-SUPERVISED LEARNING. In supervised learning, a neural network is trained on a small amount of labelled data and a large amount of unlabelled data. The labelled data helps the neural network to generalize to new examples.

GEN AI

- ↳ Subset of DL
- ↳ uses artificial neural network
- ↳ can process both labelled & unlabelled data being supervised, unsupervised and semi-supervised.
- ↳ Large Language models are a subset of AI

DEEP LEARNING MODELS

GENERATIVE

- generates new data instances based on ~~data~~ learned probability distribution of existing data
- generates new content

DISCRIMINATIVE

- type of model used to classify or predict labels for data points
- typically trained on a dataset of labelled data points.
- they learn the relationship b/w the features of the data points & the labels.
- Once trained, it can be used to predict the label for new data points.

- ↳ type of AI that creates new content based on what it has learned from existing content.
- ↳ The process of learning from existing content is called TRAINING & results in the creation of a statistical model when given a prompt.

NOTE • LLM are one type of Gen AI since they generate novel combination of text in the form of natural sounding language.

- A generative image model takes an image as input & output text in the form of natural sounding ~~text~~ language.
- A generative language model takes text as input and can output more text, an image, audio, or decisions. They are pattern matching ~~model~~ systems.

TRANSFORMER

- ↳ The power of Gen AI comes from the use of transformers.
- ↳ Transformer produced a 2018 revolution in natural language processing.
- ↳ At a high level, transformer model consists of an encoder & decoder. The encoder encodes the input sequence & passes it to the decoder, which learns how to decode the representation for a relevant tasks.

Hallucinations

- In transformer, they are words or phrases that are generated by the model that are often nonsensical or grammatically incorrect.
- Hallucinations can be caused by a no. of factors, including the model is not trained on enough data, or the model is trained on noisy or dirty data, or the model is not given enough context, or the model is not given enough constraints.
- Hallucinations can be a problem ~~for~~ for transformers because they can make the output text difficult to understand. They can also make the model more likely to generate incorrect or misleading information.

Prompt

- Short piece of text that is given to the LLM as input.
- PROMPT DESIGN is the process of creating a prompt that will generate the desired output from a LLM.

Q What are the foundation models in Gen AI?

A foundation is a large AI model pretrained on a vast quantity of data that was "designed to be adapted" (or fine tuned) to a wide range of downstream tasks, such as sentiment analysis, image captioning and object recognition.

Q What factors cause Hallucination

- (1) The model is trained on noisy or dirty data
- (2) The model is not given enough context
- (3) The model is not trained on ~~enough~~ enough data

GENERATIVE AI FOR EVERYONE

Gen AI

AI that can produce high quality content, such as text, images & audio

Q Accurate description of LLM

It generates text by repeatedly predicting the next word.

• Because LLM can hallucinate (make up facts) it is best to fact-check the response from an LLM before using it in situations where factual accuracy is important.

• Gen AI refers to a collection of tools that can generate high quality text, images & audio, including LLM and diffusion models for image generation.

• General Purpose technologies are, by definition to be versatile & useful for wide range of tasks. This broad utility across various applications is what characterizes AI as a general purpose technology.

•

• RAG (Retrieval Augmented Generation) that gives the LLM model access external data sources.

• FINE TUNING is a technique that allows you to adapt a LLM to your task. Training the model on a specific tasks or data set

• PRETRAINING MODELS which refers to training LLM from scratch. Training the model on a large corpus of text data.

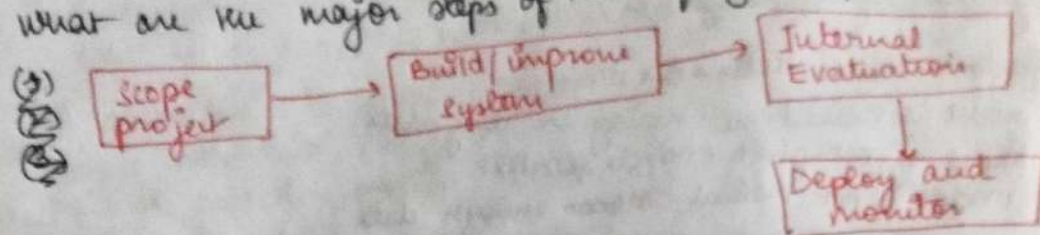
• Input - length of the prompt

• TOKEN - ~~lot~~ loosely either a word or a subpart of a word.

• LLM are very flexible, and can reason through more complex tasks if given good instructions.

- Tokens in the context of LLM refer to a unit of text. Common words are typically represented by a single token, while uncommon words may be broken into 2 or more tokens.

Q What are the major steps of the life cycle of a GenAI project



Q You are working on using LLM to summarize ~~the~~ research reports. Suppose an average report contains roughly 6000 words. Approximately how many tokens would it take an LLM to process 6,000 input words. (Assume 1 word = 1.333 tokens)

$6000 \times 1.333 = 8000 \text{ tokens}$

- RAG has 3 steps
 - (1) Given the question, it'll ~~look~~ look through a collection of documents on the benefits offered to employees that may have the answer.
 - (2) Incorporate the retrieved document of the retrieved text into an updated prompt.
 - (3) Prompt the LLM ~~with~~ with this rich prompt.
- LLM is a reasoning engine
- Development using LLM is often highly empirical, meaning experimental process

CHOOSING MODEL

(1) model size

- 1B parameter range
 - good at pattern matching
 - has basic knowledge
- 10B parameter range
 - greater knowledge of world
 - better at following instructions
- 100B+ parameter
 - better at complex reasoning
 - have rich world knowledge

(2) CLOSED VS OPEN SOURCE MODEL

• Closed source models.

- usually accessible via cloud programming interface
- relatively expensive to run because the large companies hosting these models will often have put a lot of work into serving up these API calls inexpensively.
- Downside - risk of vendor lock-in

• Open source models.

- Advantage: have full control over the model.
- lets you build the application in a way that retains full control over data privacy and data access.

RLHF

↳ Reinforcement learning from human feedback.

↳ Steps

(1) train an answer quality model.

(2) Have LLM generate a lot of answers. Further train it to generate more responses that get high scores.

↳ Reason why this called RLHF is because the scores correspond to the reinforcement or the reward that we've given the LLM for generating different answers.

NOTES

• RAG provides LLM with additional information & context from external documents that it can reason through to answer a question

Q What does the idea of using an LLM as a reasoning engine refer to?
This refers to the idea of using an LLM not as a source of information but to process information

AUGMENTATION

Help human with a task

AUTOMATION

Automatically perform a task

The potential for aug / auto a task depends on

(1) Technical feasibility - can AI do it?

(2) Business value - how valuable is it for AI to augment or automate this task?

Experiment & ~~proto~~ proto typing with web interfaces is a viable way to get started with LLM application development. This allows you to understand what is feasible before investing more time and resources in the project & team.

ARTIFICIAL GENERAL INTELLIGENCE

AGI

- ↳ Artificial General Intelligence
- ↳ AI that can do any intellectual task that a human can

DIMENSIONS OF RESPONSIBLE AI

- FAIRNESS Ensuring AI does not perpetuate or amplify biases.
- TRANSPARENCY Making AI systems and their decisions understandable.
- TO STAKEHOLDERS IMPACTED
- PRIVACY Protecting user data & ensure confidentiality
- ~~Security~~ ^{SECURITY} Safeguard AI system from malicious attacks.
- ETHICAL USE Ensuring AI is used for beneficial purposes.

NOTE

- RLHF trains model to produce output that better aligns with human preferences, including honesty, helpfulness, & harmlessness. The process can reduce biases in a LLM output.
- AI automates task not jobs

GEN AI FOR EVERYONE

FOUNDATION MODEL

- ↳ Foundation Model is like a virtual hand that has already learned a lot about language, image or even code before it starts working on a specific task.
- ↳ Foundation Models provide a starting point for building more specialized AI models.

NOTES

Q What is the primary goal of Gen AI?
To generate new & original data.

- GANs Generative ~~Adversarial~~ Adversarial Networks
- VAEs Variational Auto encoders

Q How does Gen AI impact organizational efficiency?
By automating repetitive tasks.

- Q Which of the following is a key consideration for implementing aspects of Gen AI
Ensuring unbiased model outputs.
- Q Which of the following are important considerations for ethical deployment and responsible practices in Gen AI projects?
(1) Transparency & explainability of AI system
(2) Bias & fairness in AI algorithm.
- Q What distinguishes Gen AI from other types of AI algorithms, such as discriminative algorithms?
Gen AI focuses on generating new data, while discriminative algorithms focus on classifying existing data.
- Q Assess the impact of Gen AI on organizational efficiency
(1) By optimizing resource allocation and streamlining complex business operations
(2) By speeding up the development and deployment of AI applications through automated model generation
(3) By automating decision-making processes & reducing the need for human intervention
(4) By generating realistic synthetic data for training AI models & reducing the reliance on large, labeled datasets
- Q
- GPT can generate content that lacks coherence and structure
 - ~~GPT~~ GPT — Generative Pre-trained Transformer
- Q What is an advantage of using AI technology
Increased efficiency and productivity
- Q Why is AI customized for different use cases or needs.
To optimize AI's performance & adaptability
- Q Which of the following is a key feature of coherent generation using Gen AI models like GPT-3.
The ability to generate diverse and contextually relevant text based on given prompts
- Q How can prompt engineering enhance content generation using Gen AI models like GPT-3.
By crafting specific & contextually rich prompts to influence the quality and relevance of the generated content

- Q How does GPT 4 work?
GPT utilizes transformer model & large-scale pretraining
- Q What is one potential limitation of AI when it comes to decision making based on data used
- Q What is the primary neural network architecture used in GPT?
Transformer.
- Q How is AI customized to address different use cases and needs?
By utilizing AI models & algorithms to specific tasks & requirements.
- Q How can AI be used responsibly?
By ensuring transparency, fairness, & accountability in AI systems
- Q What is one of the key steps in preparing data for training a Gen AI model?
Gathering and preprocessing the data.

CREATING A CHATBOT

- (1) Defining a Objective
- (2) Data gathering and Preparation
- (3) Language Model selection
- (4) Training the Chatbot
- (5) Integrating the Chatbot
- (6) Testing and Refinement

BUILD A CHATBOT USING GENAI

- (1) Define the Chatbot's Purpose & Scope
 - Clearly identify the goals & objective of your chatbot
 - Determine the target audience & what tasks the chatbot will handle.
- (2) Choose a Chatbot Platform or Framework
 - Research and select a user-friendly chatbot development platform.
 - Look for platform that offer pre-built templates & easy customization options.

- (3) Design Conversation
 - Plan out different user interaction & conversation.
 - Create a flowchart or diagram to visualize the chatbot's responses based on user inputs.
- (4) Teach the Chatbot to Understand
 - Use the platform's interface to define user intents (what user wants) and entities (specific information).
 - Train the chatbot by providing sample conversation & teaching it how to respond to different inputs.
- (5) Integrate with Backend Service.
 - Connect your chatbot to relevant databases or APIs to fetch information.
 - If needed, seek assistance from developers or use platform features for integration.
- (6) Test and Improve
 - Test your chatbot by having conversation & checking its response.
 - Refine and update your chatbot based on user feedback & common questions.
- (7) Deploy and monitor.
 - Launch your chatbot on your website, social media platforms, or messaging apps.
 - Monitor its performance & gather feedback to make further improvements.

ETHICAL CONSIDERATION

- (1) Data Privacy & Security
- (2) Bias & Fairness
- (3) Responsible Use
- (4) Transparency & Explainability
- (5) Ethical Framework & Guidelines
- (6) Regular Monitoring & Evaluation
- (7) User Consent & Feedback.

IMPORTANT ASPECTS

- (1) Bias Identification
- (2) Fairness Metrics
- (3) Bias Mitigation
- (4) User Feedback & Iterative Improvement
- (5) Diverse & Inclusive Training Data
- (6) External Review & Auditing
- (7) Ongoing Monitoring & Evaluation

GENERATIVE AI PRIMER

- Augmented Intelligence (AI+) is about augmenting & simplifying human creativity & problem solving.

8 How do we ACHIEVE AUGMENTED INTELLIGENCE (AI+)

- A Aid Human Coordination
- C ~~cut~~ Cut out tedious tasks
- H Help provide a safety net to make sure things aren't missed
- IEV Inspire better problem solving and creativity
- E Enable great ideas to scale faster

INTRODUCTION TO LLM

LLM

- ↳ large, general-purpose language models that can be pre-trained and then fine-tuned for specific purposes.

NOTE

- PaLM (Pathways Language Model)
- LaMDA (Language Model for Dialogue Applications)

PROMPT DESIGN

Prompt involve instruction & context passed to a language model to achieve a desired task.

PROMPT ENGINEERING

The practice of developing & optimizing prompts to efficiently use language models for a variety of applications

- There are 3 main kinds of LLM, each needs prompting in a different way

(1) GENERIC (OR Raw) LANGUAGE MODELS

These predict the next word (technically token) based on the language in the training data

(2) INSTRUCTION TUNED

Trained to predict a response to instructions given in the input

(3) DIALOG TUNED

Trained to have a dialog by predicting the next response

TUNING

- ↳ The process of adapting a model to a new domain or set of custom use cases by training the model on new data

FINE TUNING

- ↳ Using your own dataset of and retain the model by tuning every weight in the LLM. This requires a big training job (like really big) & hosting your own fine-tuned model.

PETM

- ↳ Parameter Efficient Tuning Methods
- ↳ methods for tuning a large language model on your own custom data without duplicating the model.

NOTES

- LLM
type of AI that can generate human-quality text. LLMs are trained on massive datasets of text code, and they can be used for many tasks, such as writing, translating and coding.

Q What are LLM benefits?

- (1) They can generate human-quality text
- (2) They can be used for a variety of task.
- (3) They can be trained on massive datasets of text & code.
- (4) They are constantly improving

Q What are some of the challenges of using LLMs?

- (1) They can be used to generate harmful content
- (2) They can be expensive to train
- (3) They can be ~~be~~ biased.

PROMPT ENGINEERING FOR CHATGPT: MODULE 1