



Generative AI course

1. What is Generative AI?
2. Roadmap of Generative AI?
3. What things do we need to learn in a Generative AI course?
4. What are the foundational models in Generative AI?
5. What is the future scope of Generative AI?
6. What is the eligibility for a Generative AI course?
7. Why should you learn Generative AI?
8. What are the types of Generative AI compared to other AI types?
9. What are the advantages and disadvantages of Generative AI?

what is Generative AI

Generative AI refers to AI technologies that can create new content, ideas, or data that are coherent and plausible, often resembling human-generated outputs. It has a plethora of practical applications in different domains such as computer vision, natural language processing, and music generation.

Roadmap for Generative AI

1. Foundation in Machine Learning and Deep Learning

- Ensure proficiency in probability, statistics, linear algebra, and calculus.
- Gain hands-on experience with programming languages like Python/R.
- Familiarize yourself with supervised and unsupervised learning algorithms.
- Build machine learning models on tabular datasets.

2. Deep Learning Mastery

- Develop a solid understanding of deep learning architectures such as MLPs, RNNs, LSTMs, GRUs, and CNNs.
- Gain proficiency in at least one deep learning framework like Keras, TensorFlow, PyTorch, or FastAPI.
- Train MLPs on tabular datasets.
- Construct RNNs and CNNs for unstructured data (text and image).

3. Advanced Deep Learning Techniques

- Learn about pretrained models for image data and their types.
- Understand language models and build them using LSTMs/GRUs.
- Explore attention mechanisms and their applications.
- Study autoencoders and GANs architectures and train these models on datasets.

4. Generative Models for NLP:

- Discover Large Language Models (LLMs) like Transformers, BERT, GPT, PaLM, etc.

- Understand how to use LLMs for downstream tasks: finetuning, zero-shot, one-shot, and few-shot learning.
- Learn best practices for training LLMs, including scaling laws and efficient training mechanisms.
- Explore techniques to pre-train LLMs on domain-specific data.
- Implement different techniques to finetune LLMs for downstream tasks.
- Study optimization techniques like Adapters, LoRA, QLoRA, etc., to accelerate finetuning.
- Understand deployment considerations (LLMops) for deploying LLMs in production.
- Explore cutting-edge models like ChatGPT and BARD and understand their training process, including reinforcement learning from human feedback, supervised fine-tuning, and prompt engineering.
- Gain hands-on experience with LLM frameworks like LangChain, AutoGPT, Vector DB etc.

Disadvantages :

1. **Source Ambiguity:** Generative AI may not always attribute content to its original source, leading to ambiguity in authorship and ownership.
2. **Bias Assessment Difficulty:** Assessing bias in original sources can pose challenges, complicating efforts to ensure fairness and accuracy in generated content.
3. **Realism Challenges:** The lifelike quality of generated content can obscure inaccuracies, making it more difficult to distinguish between genuine and

fabricated information.

4. **Tuning Complexity:** Adapting generative AI for new contexts and scenarios can be intricate, requiring nuanced adjustments to ensure optimal performance.
5. **Potential for Oversights:** Despite advancements, there remains a risk that generative AI results may inadvertently perpetuate biases, prejudices, and hateful ideologies, warranting careful scrutiny and mitigation measures .

What are the benefits of generative AI?

- Automating the manual process of writing content.
- Reducing the effort of responding to emails.
- Improving the response to specific technical queries.
- Creating realistic representations of people.
- Summarizing complex information into a coherent narrative.
- Simplifying the process of creating content in a particular style.

What are the concerns surrounding generative AI?

- It can provide inaccurate and misleading information.
- It is more difficult to trust without knowing the source and provenance of information.
- It can promote new kinds of plagiarism that ignore the rights of content creators and artists of original content.
- It might disrupt existing business models built around search engine optimization and advertising.

- It makes it easier to generate fake news.
- It makes it easier to claim that real photographic evidence of a wrongdoing was just an AI-generated fake.
- It could impersonate people for more effective social engineering cyber attacks.

What are some examples of generative AI tools?

- **Text generation tools** include **GPT**, **Jasper**, **AI-Writer**, and **Lex**.
- **Image generation tools** include **DALL-E 2**, **Midjourney**, and **Stable Diffusion**.
- **Music generation tools** include **Amper**, **Dadabots**, and **MuseNet**.
- **Code generation tools** include **CodeStarter**, **Codex**, **GitHub Copilot**, and **Tabnine**.
- **Voice synthesis tools** include **Descript**, **Listnr**, and **Podcast.ai**.
- **AI chip design tool** companies include **Synopsys**, **Cadence**, **Google**, and **Nvidia**.

What are use cases for generative AI?

- Implementing chatbots for customer service and technical support.
- Deploying deepfakes for mimicking people or even specific individuals.

- Improving dubbing for movies and educational content in different languages.
- Writing email responses, dating profiles, resumes and term papers.
- Creating photorealistic art in a particular style.
- Improving product demonstration videos.
- Suggesting new drug compounds to test.
- Designing physical products and buildings.
- Optimizing new chip designs.
- Writing music in a specific style or tone.

Some of the very commonly known LLMs are:

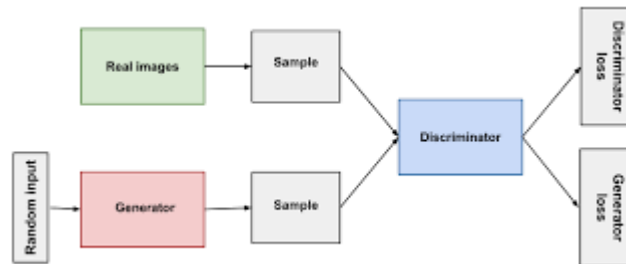
- Open AI's GPT 3, 3.5, and 4
- Google's LaMDA and PaLM
- Hugging Face's BLOOM
- Meta's LLaMA
- NVidia's NeMO LLM

Types of Models for Generative AI ?

1. **Generative Adversarial Networks (GANs):** GANs consist of two neural networks, a generator and a discriminator, trained simultaneously in a game-like scenario. The generator tries to create data (e.g., images) that are indistinguishable from real data, while the discriminator tries to

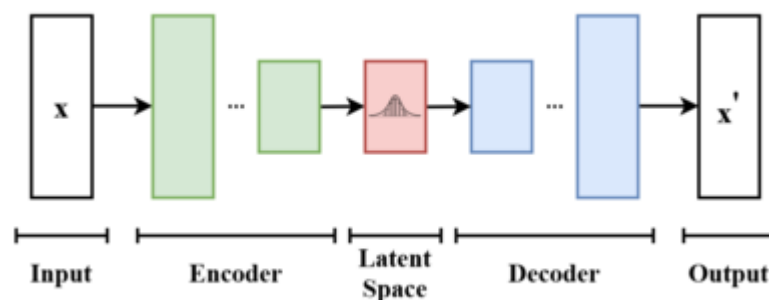
differentiate between real and generated data. This adversarial training process helps improve the quality of generated samples.

2.



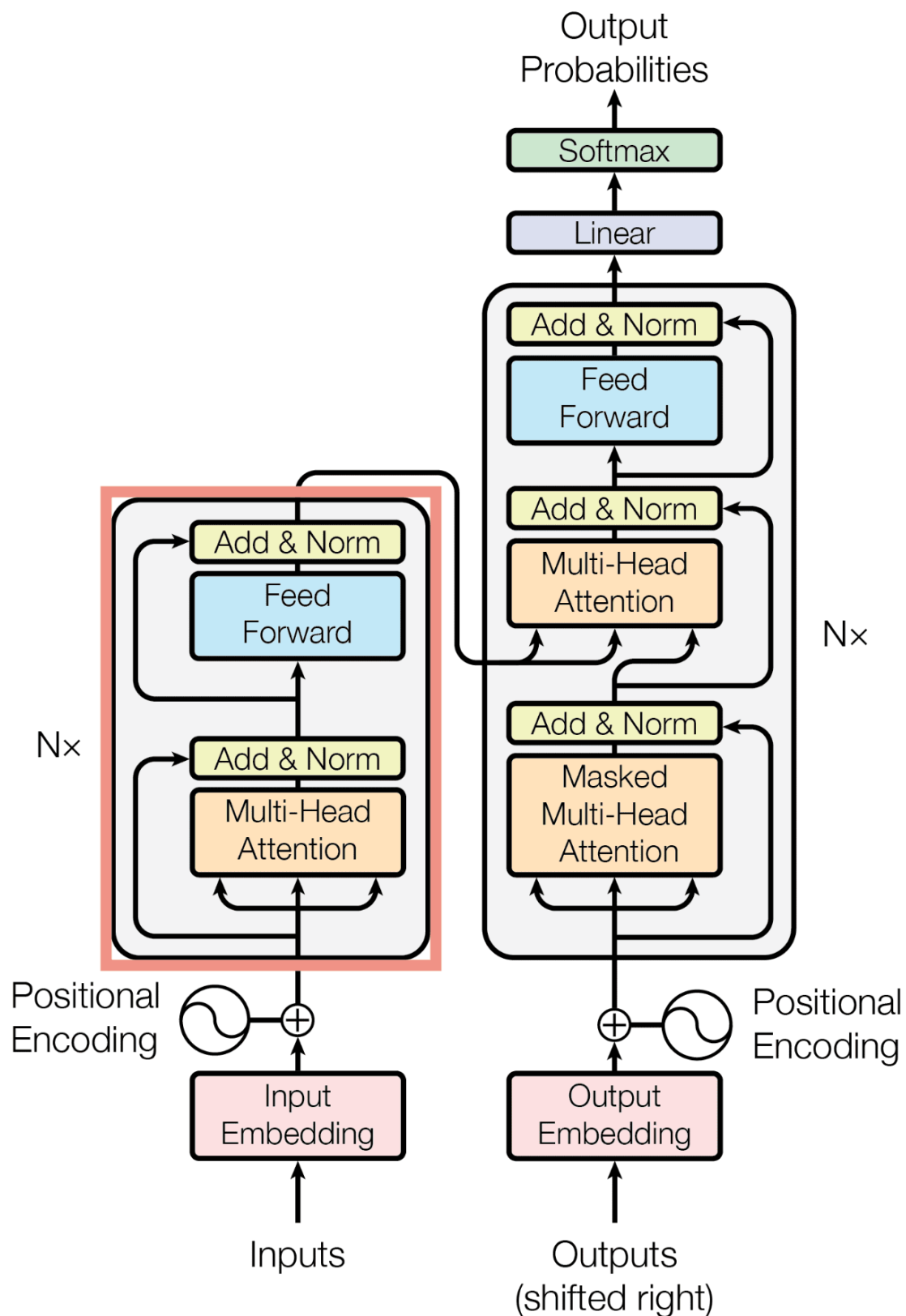
1. **Variational Autoencoders (VAEs):** VAEs are generative models based on the principles of autoencoders. They consist of an encoder network that compresses input data into a latent space representation and a decoder network that reconstructs the original data from the latent space. VAEs are trained to generate new data samples by sampling from the learned latent space distribution.

2.



1. **Transformers:** Transformers are a type of deep learning model primarily used in natural language processing (NLP). They are based on a self-attention mechanism that allows the model to weigh the importance of different words in a sequence when processing each word. Transformers have achieved state-of-the-art performance in various NLP tasks and have also been adapted for other tasks beyond NLP, such as image generation and reinforcement learning.

2.

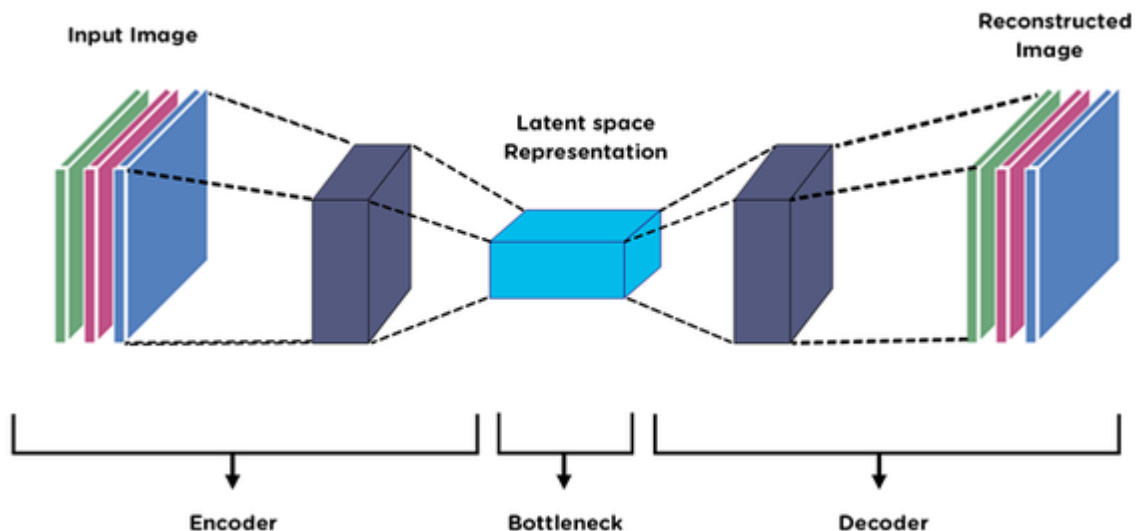


1. **Autoencoders:** Autoencoders are neural networks trained to copy their input to their output. They consist of an encoder network that compresses

the input data into a latent space representation and a decoder network that reconstructs the original input from the latent space representation.

Variants like denoising autoencoders, sparse autoencoders, and contractive autoencoders serve different purposes.

2.



Future capabilities of generative AI:

Continued Model Scaling and Efficiency Gains

- **Larger and more parameter-dense models** promise to further expand generative capabilities and general competencies.
- At the same time, **more efficient training techniques** will reduce data and computing needs to improve sustainability.
- Also, **new hardware and software architectures** could be developed to more effectively handle resource intensive workloads.

Steady Improvements in Output Quality and Accuracy

- As this technology advances, its responses may become indistinguishable from, and **possibly superior to human-generated material**.
- Improvements will be achieved through **advances in responsible development, model architecture, training techniques, data quality and diversity, and evaluation metrics and tools**.
- Also, **end users and developers contribute to improvements** by providing feedback on responses, labeling data, generating ideas, developing models, and building tools.

Multiple Layered Agent Frameworks:

- **Multiple layered agent frameworks are a complex and hierarchical structure of AI agents** designed to interact and collaborate effectively in a dynamic environment.
- This type of framework allows numerous agents to communicate with one another in order to **solve problems that are too complicated for a single agent to accurately accomplish on its own**.
- In addition, multiple AI agents working together **can dramatically increase the productivity of task completion**.
- For instance, picture a transportation network management system where sensor agents gather real-time data, data processing agents analyze it, routing agents determine optimal paths, coordination agents ensure harmonious cooperation, and supervisory agents oversee the entire operation.

Autonomous AI Agents

- One or more AI agents with **the ability to learn and adapt on their own**, without the need for human intervention.
- This will allow them to **automate tasks and make decisions in real time**.
- For example, an autonomous AI agent could be used to **monitor current events and automatically generate content on flagged topic categories**.

Increased Customization and Specialization :

- Generative AI will become more customizable, allowing for **more parameters to be fine-tuned on custom data**.
- This contextual fine-tuning improves relevance, quality and safety for specific professional use cases while retaining general advantages.
- Which will **serve as springboards for developing specialized systems** tuned on custom datasets for focused tasks, verticals and applications.
- To illustrate, a healthcare company could fine-tune a generative AI model on medical data to generate new drug candidates, predict the effectiveness of different drugs, and design new drug compounds.

Accessibility :

- Generative AI tools will **become more accessible to people with less data and expertise**.
- **Making it easier for businesses of all sizes to use this technology** by lowering barriers to access, enabling broader adoption across languages, domains, geographies and less resourced organizations.
- Specifically, the development of user-friendly apps and AI-as-a-service (ALaaS) platforms will **enable access without having to build and maintain their own infrastructure**.

Text to Everything

- The ability to **generate text, images, audio, videos, and 3D environments from a single platform**.
- This means users will **provide a single text prompt and the model will generate every content format simultaneously**.
- For example, a user could provide a prompt like "A beautiful sunset over a beach" and the model would generate a text description of the sunset, a video of the sunset, and an audio file of the sound of waves crashing on the shore.

Ultra Realistic Talking Avatars :

- Generative AI models will be able to **create videos of talking avatars that are so realistic that they are indistinguishable from real people.**
- These avatars will be **able to lip sync to speech, have natural facial expressions, subtle head movements and even generate their own dialogue.**
- **This will have a wide range of applications,** including virtual assistants, customer service representatives, and even educational characters

Impact of Generative AI on industry and jobs :

Changing Job Roles and Required Skills :

- Though AI won't replace human roles outright, it **will alter the associated required skills.**
- **workers and job requirements will need to transition into more strategic roles.**
- **This may involve upskilling to perform higher-value work** like data analysis, creative direction and strategy.

Automating Tasks to Increase Efficiency

- Conducting process analysis can be used to **identify repetitive and low strategic value tasks that are ideal for AI automation.**
- Task categories such as document review, contract analysis, customer service routines and content generation contain repetitiveness that humans find tedious.

- **free up employees for more impactful and creative responsibilities** better suited to human strengths.

Renewed Focus on Data Governance

- **Generating high-quality outputs relies on extensive training data** that is accurate, comprehensive and aligned to use cases.
- **This places greater emphasis on responsible data governance**, including ethics, integrity, monitoring, security and compliance.
- Where **data collection practices will require more comprehensive auditing** to avoid perpetuating biases.

Tight Integration into Workflows and Software

- To maximize user adoption and impact, AI systems **will integrate seamlessly into existing enterprise software, workflows and processes**.
- This integration will **emphasize on providing a high-quality user experience by actively seeking continuous feedback** to drive enhancements.
- Ultimately, **the main objective of these AI systems is to seamlessly integrate into workflows**, to feel like a natural system component rather than function as isolated tools.

Lower Barriers of Entry

- Generative AI **will make it easier for new competitors to enter the market**.
- As pre-trained models and AI-as-a-service (ALaaS) platforms become more accessible, smaller players will leverage these technologies without massive in-house investments, **somewhat leveling the playing field**.
- As a result, **this will increase market competition** and force existing enterprises to adapt quickly and differentiate themselves.

Enabling New Niche Capabilities

- With the help of generative AI, one can **scan the market for untapped areas to uncover potential niche opportunities**.

- Specifically, **to identify new market opportunities from various sources** by analyzing customer data, social media, demographics, economic trends, etc.
- By feeding all these data points together, **the AI model will be able to identify new customer needs and wants** that are not currently being met by existing products and services.
 - Then **integrate such findings into the process to generate new ideas** for products and services.

Applications Across Industries:

- **Marketing, Advertising, and Entertainment Industry:**
 1. Content Creation: Generative AI powers content creation in various forms like art, music, and literature, allowing artists and musicians to explore innovative creative directions.
 2. Video Game Development: AI-driven systems create game environments, characters, and dialogues, reducing development time and resources.
 3. Scriptwriting: Screenwriters use Generative AI to assist in scriptwriting by generating dialogues, plotlines, and character interactions.
- **Education Sector:**
 1. Personalized Learning: Generative AI adapts educational content to individual student needs by generating tailored assignments, quizzes, and study materials.
 2. Knowledge Base: AI creates exhaustive knowledge bases for students to access information in a conversational style.
 3. Virtual Labs: Generative AI powers virtual laboratories for simulating experiments and scenarios in science and engineering disciplines.
- **Healthcare Industry:**
 1. Medical Image Generation: Generative AI creates synthetic medical images for training models, enhancing diagnostic accuracy, and simulating rare medical conditions.

2. Drug Discovery: Pharmaceutical companies use Generative AI to discover new drug compounds by generating molecular structures.
 3. Personalized Medicine: AI analyzes patient data to generate personalized treatment plans, accounting for genetic factors and medical history.
- **Manufacturing Industry:**
 1. Product Design: Generative design creates optimized product designs considering materials, weight, and structural integrity.
 2. Quality Control: AI generates synthetic data for quality control testing, ensuring adherence to quality standards.
 3. Supply Chain Optimization: AI-generated forecasts help make informed decisions about production and distribution.
 - **Software & Tech Industry:**
 1. Code Generation: Generative AI assists developers by generating code snippets and templates, speeding up development.
 2. Bug Detection: AI generates synthetic test cases to identify and fix software bugs efficiently.
 3. IT Security: Generative AI models simulate cyberattack scenarios to identify vulnerabilities and enhance cybersecurity measures.

What is prompt engineering

Prompt engineering is the practice of designing inputs for AI tools that will produce optimal output

why we need to learn prompt engineering ?

Enhanced accuracy and relevance

Improved decision making

Personalized customer experiences

Efficient resource utilization

Ethical considerations and bias mitigation

There are 8 prompt engineering methods:

(1) Zero-Shot Learning, (2) One-Shot Learning, (3) Few-Shot Learning, (4) Chain-of-Thought Prompting, (5) Iterative Prompting, (6) Negative Prompting, (7) Hybrid Prompting, and (8) Prompt Chaining .

1.

Zero-Shot Learning: This involves giving the AI a task without any prior examples. You describe what you want in detail, assuming the AI has no prior knowledge of the task.

- . **One-Shot Learning:** You provide one example along with your prompt. This helps the AI understand the context or format you're expecting.
3. **Few-Shot Learning:** This involves providing a few examples (usually 2–5) to help the AI understand the pattern or style of the response you're looking for.
4. **Chain-of-Thought Prompting:** Here, you ask the AI to detail its thought process step-by-step. This is particularly useful for complex reasoning tasks.
5. **Iterative Prompting:** This is a process where you refine your prompt based on the outputs you get, slowly guiding the AI to the desired answer or style of answer.
6. **Negative Prompting:** In this method, you tell the AI what not to do. For instance, you might specify that you don't want a certain type of content in the response.
7. **Hybrid Prompting:** Combining different methods, like few-shot with chain-of-thought, to get more precise or creative outputs.
8. **Prompt Chaining:** Breaking down a complex task into smaller prompts and then chaining the outputs together to form a final response.

Key elements of a prompt :

- **Instruction.** This is the core directive of the prompt. It tells the model what you want it to do. For example, "Summarize the following text" provides a clear action for the model.
- **Context.** Context provides additional information that helps the model understand the broader scenario or background. For instance, "Considering the economic downturn, provide investment advice" gives the model a backdrop against which to frame its response.
- **Input data.** This is the specific information or data you want the model to process. It could be a paragraph, a set of numbers, or even a single word.
- **Output indicator.** Especially useful in role-playing scenarios, this element guides the model on the format or type of response desired. For instance, "In the style of Shakespeare, rewrite the following sentence" gives the model a stylistic direction.

Technical skills for prompt engineering :

- **Understanding of NLP.** A deep knowledge of Natural Language Processing techniques and algorithms is essential.
- **Familiarity with LLMs.** Experience with models like GPT, PaLM2, and other emerging models their underlying architectures.
- **Experimentation and iteration.** Ability to test, refine, and optimize prompts based on model outputs.

Non-technical skills for prompt engineering

- **Communication.** The ability to convey ideas, collaborate with teams, and understand user needs.
- **Subject Matter Expertise.** Depending on the application, domain-specific knowledge can be invaluable.
- **Language Proficiency.** Mastery over language, grammar, and semantics to craft effective prompts.
- **Critical Thinking.** Evaluating model outputs, identifying biases, and ensuring ethical AI practices.
- **Creativity.** Thinking outside the box, experimenting with new prompt styles, and innovating solutions.

Elements of a Prompt

A prompt can include different elements:

- **Instruction** — It's a specific task or direction for the model to follow. :
- **Context**— It provides extra information or context to help the model generate better responses.
- **Input Data** — It refers to the question or input for which we want the model to provide a response.
- **Output Indicator** — It indicates the desired type or format of the output.

Prompt Engineering use cases:

1. Text summarization
2. Information Extraction
3. Question Answering
4. Text Classification
5. Code generation
6. Chatbots and Virtual Assistants

Effective prompt engineering involves several best practices:

1. **Clearly Define Desired Responses:** Specify the scope of the desired response to prevent misinterpretations by LLMs.
2. **Be Specific and Explicit:** Avoid vagueness; provide clear, explicit instructions to guide LLMs effectively.
3. **Balance Simplicity and Complexity:** Find the right balance between simple and complex prompts to ensure clarity without overwhelming Large Language Models.
4. **Iterate and Experiment:** Prompt engineering is iterative; test different prompts and refine them based on results.
5. **Train and Evaluate:** Consider additional training to improve AI model performance if necessary, creating custom models for specific tasks.

1.Task		Critical
2.Context		Important
3.Examples		
4.Persona		
5.Format		Added Value
6.Tone		

Task

"Task" is the actual prompt, but without the constraints.

Context

"Context" is the actual detailed content and it is the 2nd most important and essential part of your prompt

Examples

"Examples" is also an important part of prompting since giving clear and concise examples tells the large language model about the desired outcome and style

Persona

"Persona" is the moment when you say: "Write Like You Are" or "Create in the style of" to create the content as if that person is writing it.

Format

"Format" is about in what shape or form you want your results to be.

Tone

"Tone" is when you say: professional, friendly, casual, concise, emotional, etc. It is the feeling that you want to convey to your readers.

Prompting Principle :

- **Principle 1: Write clear and specific instructions**

Tactics

1: Use delimiters to clearly indicate distinct parts of the input

- Delimiters can be anything like: ```, """, < >, `<tag> </tag>`, `:`
 - 2. Ask for the structured output
 - 3. ask model to check whether condition satisfied or not
 - 4. few shot prompting
- **Principle 2: Give the model time to "think":**
 - 1: Specify the steps required to complete a task
 - 2. : Instruct the model to work out its own solution before rushing to a conclusion

Model Limitations: Hallucinations

How to Iterate your prompt for getting the actual output from it :

1. The text is too long
2. Text focuses on the wrong details
3. description is not as you wanted

2. Summarizing :

1. Summarize with a word/sentence/character limit

2. summarize with your focus point like price , values

3. also you can summarize list of review by prompting

3. Inferring :

1. Entity Recognition: Identifying and categorizing entities mentioned in text, such as people, organizations, locations, dates, and more.
2. Part-of-Speech Tagging: Assigning grammatical categories (e.g., noun, verb, adjective) to words in a sentence.
3. Named Entity Recognition (NER): Identifying proper nouns that refer to specific entities, such as names of people, places, organizations, etc.
4. Coreference Resolution: Determining which words or phrases in a text refer to the same entity.
5. Semantic Role Labeling: Identifying the roles of words or phrases in relation to a predicate (e.g., who did what to whom).
6. Textual Entailment: Determining the logical relationship between two pieces of text, such as whether one sentence entails or contradicts another.
7. Question Answering: Generating accurate responses to questions posed in natural language based on a given context.

sentiment analysis

identify the emotion

extract multiple things

4. Transforming :

language correction
language translation
Tone transformation
format conversion
spell and grammar checker

5. Expanding : convert short information into longer one:

replying anyone
writing essay
explaining topic

6. Build your Bot :

Generative ai Terms and Terminology :

LIFE CYCLE OF GENERATIVE AI PRODUCT

. Objective: Define Business Use Case

In this initial phase, the primary focus is on articulating the specific business challenge where leveraging a Large Language Model (LLM) strategically for enterprise applications is crucial. Key selection criteria include:

- Identifying a suitable use case that aligns with business objectives.
- Ensuring the capabilities of the chosen LLM align with the identified challenge.
- Assessing execution feasibility in terms of available resources and required skills.
- Hypothesizing potential business value generated from implementing the LLM solution.
-

2. Select Model: Choose the Right LLM Model for the Business Case

This stage involves strategically selecting the most suitable LLM model tailored to meet the unique demands of the identified business use case.

Considerations include:

- Choosing between models such as GPT-3 or LLMA for text generation, based on specific requirements.
- Evaluating models specialized in tasks like Sentiment Analysis or multilingual translation, such as M2M-100.
- Aligning the choice with the production or research environment and scalability needs.
-

3. Adapt and Harmonize: Prompt Engineering

The focus here is on crafting input prompts strategically to guide LLMs in generating precise and relevant outputs. Benefits include:

- Generating accurate and contextually relevant responses to queries.
- Providing contextual guidance to mitigate biases and improve response quality.
- Saving time and resources by streamlining LLM output.
-

4. Adapt and Harmonize: Contextual Enrichment with RAG

Introducing Retrieval Augmented Generation (RAG) enriches LLM performance by providing additional context during text generation, leading to:

- Enhanced accuracy, depth, and relevance of generated outputs.
- Reduction of inaccuracies or hallucinations in responses.
- Improved overall LLM performance.
-

5. Adapt and Harmonize: Model Refinement through Instruction Fine-tuning

Instruction fine-tuning involves supervised learning to enhance LLM performance across different tasks by:

- Training the model with labeled examples to improve task-specific performance.
- Utilizing techniques like LoRA fine-tuning or soft prompts to adjust model behavior.
- Enhancing model adaptability and performance through iterative refinement.
-

6. Adapt and Harmonize: User-Centric Improvement via Feedback Integration

Iteratively incorporating user feedback ensures that LLMs align with user expectations and ethical standards by:

- Establishing a continuous improvement cycle based on user feedback.
- Addressing biases and privacy concerns to enhance user trust.
- Ensuring LLM development aligns with user requirements and ethical guidelines.
-

7. Adapt and Harmonize: Comprehensive Model Assessment with LLM Performance Evaluation

Assessing LLM performance involves employing specialized metrics like ROUGE and BLEU to:

- Evaluate generated outputs in tasks such as summarization and translation.
- Gain nuanced insights into model efficacy while acknowledging evaluation complexities.

- Utilize multiple benchmarks like GLUE, SuperGLUE, MMLU, BIG-Bench, and HELM for comprehensive assessment.
-

8. Application Integration: Optimize and Deploy LLMs

Transitioning from training to real-world application involves optimizing LLMs for performance and resource efficiency through techniques such as:

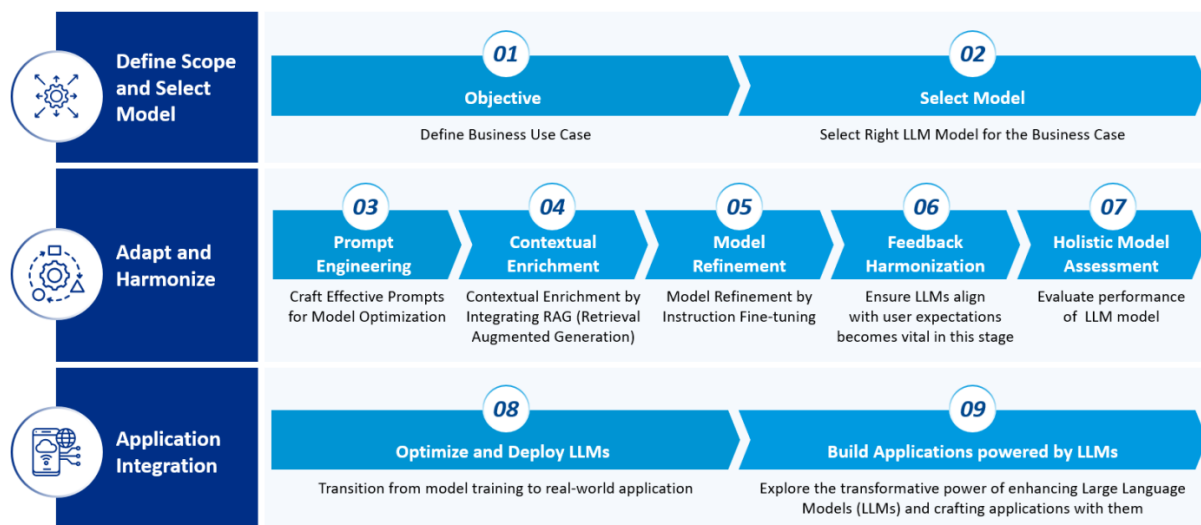
- Model distillation to train smaller models mimicking larger ones.
- Quantization to reduce memory usage and improve efficiency.
- Pruning to streamline the model without sacrificing accuracy.

Deploying optimized LLMs can occur on-site, on the cloud, or on edge devices as per requirements.

9. Application Integration: Develop Applications Empowered by LLMs

In this final stage, applications are tailored and refined using LLMs to ensure precise responses by:

- Integrating new data and considering domain-specific contexts.
- Offering accurate answers for tasks ranging from chatbots to decision-making systems.
- Fine-tuning models based on user input to enhance performance and relevance.



1. Data Preparation: Data Preprocessing and Filtering

In this stage, the focus is on preparing the data for training by performing preprocessing and filtering. This involves:

- Cleaning and organizing the data to remove noise and inconsistencies.
- Preprocessing text data by tokenizing, normalizing, and removing stop words or irrelevant information.
- Filtering the data to ensure it meets quality standards and is relevant to the model's objectives.

2. Train Model: Select a Foundational Model and Begin Training

Here, the foundational model is selected based on the task requirements, and training is initiated by providing the model with tokenized data. Steps include:

- Choosing an appropriate foundational model, considering factors like architecture, size, and pre-trained weights.
- Providing tokenized input data to the model to start the training process.
- Monitoring and adjusting training parameters to optimize model performance and convergence.

3. Validate the Model: Create Model Card

Validation involves creating a model card to provide comprehensive information about the trained model. Key components include:

- Documenting model architecture, parameters, and training data.
- Describing model performance metrics and evaluation methods.
- Disclosing any biases, limitations, or ethical considerations associated with the model.

4. Tuning the Model: Fine-Tuning for Optimization

Fine-tuning the model aims to optimize its performance further based on validation results. Steps include:

- Adjusting model parameters, such as learning rate or batch size, based on validation feedback.
- Fine-tuning the model on specific tasks or domains to improve accuracy and generalization.
- Iteratively refining the model to achieve desired performance levels.

5. Deploy the Model:

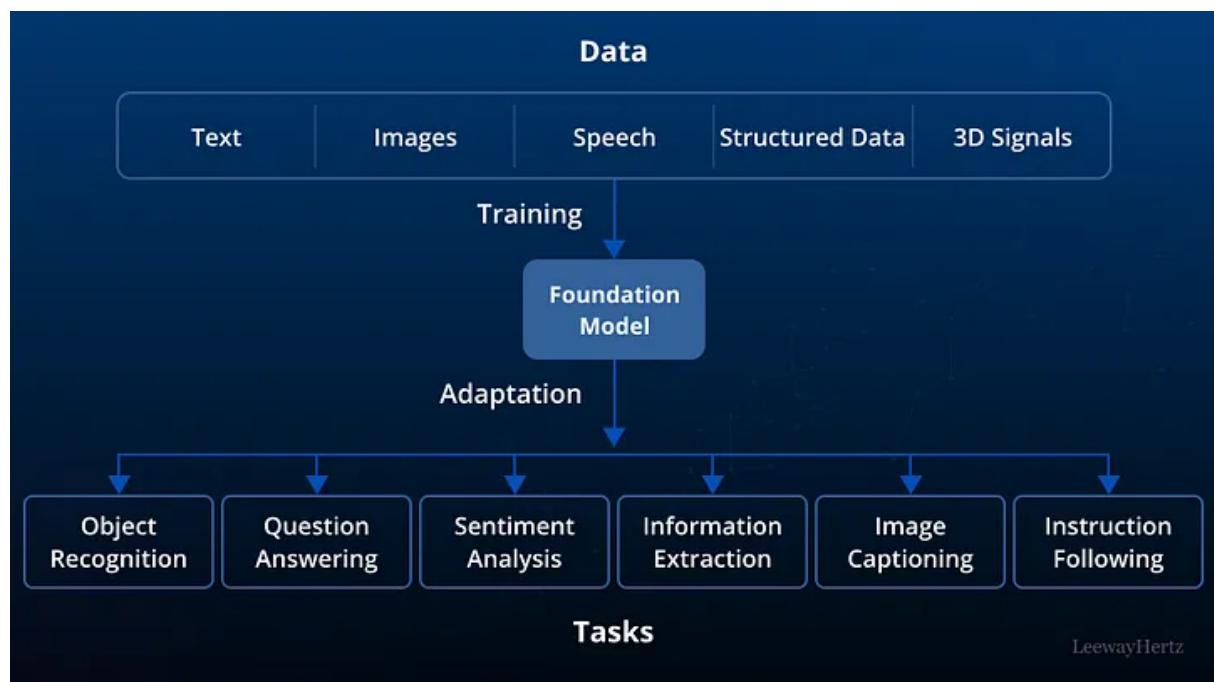
Deployment involves preparing the model for real-world usage and ensuring its readiness for deployment. Actions include:

- Reviewing and correcting any grammatical errors or inconsistencies in the model documentation.
- Highlighting key stages or actions in the deployment process by formatting them in bold for clarity.
- Deploying the model to production environments or platforms for end-user access and utilization.

What Defines a Foundation Model?

Foundation models (FMs) stand as monumental achievements in the field of machine learning (ML), reshaping the landscape for data scientists. Unlike traditional approaches that require building AI systems from scratch, FMs provide a starting point, accelerating the development of new ML applications. Coined by researchers, the term "foundation model" refers to large neural networks trained on extensive datasets. These models excel in a myriad of

tasks, including language understanding, text generation, image processing, and natural language conversation.



The Unique Traits of Foundation Models

Foundation models possess an unparalleled adaptability, allowing them to excel across diverse tasks with remarkable accuracy, driven by input prompts. Unlike conventional ML models confined to specific tasks, foundation models boast a broad spectrum of capabilities, from natural language processing (NLP) to image classification. Their size and versatility set them apart, enabling them to tackle complex challenges with finesse.

Applications and Use Cases

The versatility of foundation models lends itself to a multitude of applications, including:

- Enhancing customer support experiences
- Facilitating language translation services
- Generating diverse content
- Streamlining copywriting tasks
- Enabling accurate image classification
- Empowering high-resolution image editing

- Extracting insights from documents
- Advancing robotics and automation
- Transforming healthcare practices
- Revolutionizing autonomous vehicle technologies

Challenges in Utilizing Foundation Models

Despite their prowess, foundation models present several challenges:

- High infrastructure requirements and resource-intensive training processes.
- Complex integration into software stacks, necessitating expertise in prompt engineering and fine-tuning.
- Limited contextual comprehension, leading to occasional inaccuracies in responses.
- Potential for unreliable or inappropriate answers, especially on unfamiliar subjects.
- Risks of bias and exposure to inappropriate content due to training data characteristics.

Categories of Foundation Models

Foundation models encompass a variety of types, broadly classified into:

- **Language Models:** Specialized in processing and understanding natural language, enabling tasks like translation and text generation. Examples include BERT, GPT-3, and T5.
- **Computer Vision Models:** Geared towards analyzing visual data, performing tasks such as image classification and object detection. Notable models include ResNet, VGG, and Inception.
- **Multimodal Models:** Designed to process both text and visual data, facilitating tasks like image captioning and visual question answering. Examples include DALL-E 2, Flamingo, and Florence.

Examples of Foundation Models

pre training :

- Data Collection
 - General Text Data
 - Specialized Data
- Data Preprocessing
 - Quality Filtering
 - Deduplication
 - Tokenization
- How Does Pretraining Affect LLMs?
 - Mixture of Sources
 - Amount of Pretraining Data
 - Quality of Pretraining Data
- Architecture for Pretraining
 - Encoder Decoder Architecture
 - Causal Decoder Architecture
 - Prefix Decoder Architecture
 - Emergent Architectures
- Some Notes on Configurations
 - Layer Normalization
 - Attention
 - Positional Encoding
- Pretraining Tasks
 - Language Modeling
 - Denoising Autoencoding
 - Mixture-of-Denoisers

Pre-training is the foundation for the power of Large Language Models (LLMs) . By training on lots of text data, these models get good at understanding and generating language.

- **Importance of Pre-training Data:**

- LLMs need to be trained on a lot of data to become really capable.
- The size and quality of this data are super important. Good data helps the model achieve better capabilities.

There are two main types of data used to train Large Language Models (LLMs): general data and specialized data.

- **General Data:**

- Examples: Web pages, books, chat logs.
- Why it's used: It's widely available, varied, and helps LLMs get good at general language understanding and adaptability.

- **Specialized Data:**

- Examples: Datasets in multiple languages, scientific content, and programming code.
- Purpose: Helps LLMs become experts in specific areas or tasks.

General Text Data:

Most Large Language Models (LLMs) use general-purpose data. Let's look at three key types:

- **Webpages:**
 - What it offers: A wide range of data from the internet that provides diverse language knowledge.
 - Example Dataset: CommonCrawl.
 - Issues: Web data can have both high-quality text (like Wikipedia) and low-quality text (like spam). It's crucial to clean and process this data to ensure quality.
- **Conversation Text:**

- Why it's useful: Helps LLMs get better at conversation and question-answering.
- Example Dataset: PushShift.io Reddit corpus.
- How it's used: Conversations from online platforms are structured into tree-like formats to capture responses. This allows multi-party chats to be split into smaller conversations for training.
- Challenges: Relying too much on dialogue data can lead to problems. For instance, LLMs might misunderstand certain phrases as conversation starters, affecting their response quality.
- Books:
 - Why they matter: Books offer formal and lengthy texts which help LLMs understand complex language structures, grasp long-term context, and produce coherent narratives.
 - Example Datasets: Books3 and Bookcorpus2 found in the Pile dataset. In essence, these general data sources help LLMs understand and generate varied and natural language, but each source comes with its unique strengths and challenges.

Specialized Text Data:

Specialized data helps LLMs get better at certain specific tasks. Here are three types of specialized data:

- Multilingual Text:
 - Purpose: Improving language understanding and generation across multiple languages.
 - Example Datasets: BLOOM (covers 46 languages) and PaLM (covers 122 languages).
 - Benefit: These models are great at tasks like translation, multilingual summaries, and multilingual Q&A, sometimes even outperforming models trained just on target language data.
- Scientific Text:
 - Why it's used: Helps LLMs understand scientific knowledge.

- Source: Materials like arXiv papers, scientific textbooks, math websites, and more.
- Challenges & Solutions: Scientific texts have things like math symbols and protein sequences. To handle this, they're specially tokenized and pre-processed to fit a format LLMs can use.
- Benefits: LLMs trained on scientific texts are better at scientific tasks and reasoning.
- Code:
 - Why it's relevant: Training LLMs on code helps with program synthesis, a popular research area.
 - Current State: Even powerful LLMs, like GPT-J, find it tough to produce good, accurate programs.
 - Source: Code can come from Q&A communities like Stack Exchange or public software repositories like GitHub. This includes actual code, comments, and documentation.
 - Challenges & Solutions: Code has its own syntax and logic, and is very different from regular text. Training on code, though, might give LLMs complex reasoning skills.
 - Benefits: When tasks are formatted like code, LLMs can produce more accurate answers. In short, using specialized data gives LLMs specific skills, from understanding multiple languages to generating code.

Filtering :

Tokenization	Splitting the raw text into individual tokens or subword units, which can be fed into the model as input. This can be done using various algorithms, such as whitespace-based, rule-based, or statistical methods.
Quality Filtering	Removing low-quality data from the corpus using either classifier-based or heuristic-based approaches. Classifier-based methods train a binary classifier with well-curated data as positive instances and sample candidate data as negative instances, and predict the score that measures the quality of each data example.

	Heuristic-based methods use rules or heuristics to filter out low-quality data based on certain criteria.		
Deduplication	Duplicate documents can arise from various sources, such as web scraping, data merging, or data augmentation, and can lead to several issues, such as overfitting, bias, or inefficiency. To address these issues, existing studies mostly rely on the overlap ratio of surface features (e.g., words and n-grams overlap) between documents to detect and remove duplicate documents containing similar contents. Furthermore, to avoid the dataset contamination problem, it is also crucial to prevent the overlap between the training and evaluation sets, by removing the possible duplicate texts from the training set. It has been shown that the three levels of deduplication (i.e., document-level, sentence-level, and token-level) are useful to improve the training of LLMs, which should be jointly used in practice.		

How does pretraining affect LLMs?

Effect of pretraining on LLMs can be broadly categorized into three categories

- Mixture of sources
- Amount of pretraining data

- Quality of pretraining data

Pretraining Tasks

There are two commonly used pretraining tasks:

1. Language Modelling
2. Denoising Autoencoding

Language Modeling in LLMs:

Language Modeling (LM) is a foundational task in natural language processing. It entails predicting the next token in a sequence based on its history. When applied to Large Language Models (LLMs), especially decoder-only models, it serves as a pivotal pre-training task

Denoising Tasks in LLM Pre-training:

The discussion revolves around denoising tasks, which are prominent pre-training objectives for Large Language Models (LLMs). Let's break down the content:

1. Denoising Autoencoding (DAE):

- **Definition:** In the Denoising Autoencoding task, parts of the input text are intentionally corrupted by replacing certain spans. The objective is to train the model to recover the original, uncorrupted tokens

Influence of Architecture and Pre-training Tasks on LLMs:

1. Architecture Choice:

- **Discussion:** Early literature on pre-trained language models extensively discussed architectural effects. However, many LLMs use the causal decoder architecture, with limited theoretical analysis on its advantages.
- **Causal Decoder and LM Objective:**

- LLMs using a causal decoder architecture with a language modeling (LM) objective have shown strong zero-shot and few-shot generalization capabilities.
- Without multi-task fine-tuning, the causal decoder performs better in zero-shot scenarios than other architectures.
- GPT-3, a popular model, confirmed that large causal decoders can be effective few-shot learners.
- Techniques like instruction tuning and alignment tuning can enhance the performance of large causal decoder models.
- **Scaling Law:**
 - Causal decoders benefit from scaling laws: increasing model size, dataset size, and computation can notably improve performance.
 - In-depth studies on encoder-decoder models, especially at larger scales, are needed.
- **Future Research:**
 - More research is necessary to understand how architecture and pre-training task choices affect LLM capacity. Particular interest is in encoder-decoder architectures. Additionally, detailed LLM configurations deserve more attention.

1. Long Context:

- **Context Limitation:** Transformers have traditionally been constrained by context length due to quadratic computational costs in terms of time and memory.
- **Growing Demand:** With increasing needs for long context windows in tasks like PDF processing and story writing, models are evolving. For instance, ChatGPT has expanded its context window from 4K tokens to 16K tokens, and GPT-4 has been extended to 32K tokens.
- **Extrapolation:**
 - This refers to an LLM's ability to handle input texts that are longer than the maximum length seen during training.

- Position embedding techniques like RoPE and T5 bias have displayed extrapolation abilities. For example, LMs equipped with ALiBi have demonstrated consistent performance on sequences much longer than training sequences. Furthermore, the xPos method seeks to enhance the extrapolation capability of RoPE.
- **Efficiency:**
 - To address the quadratic computational challenge, various studies have proposed more efficient attention computation methods, such as sparse or linear attentions.
 - FlashAttention improves efficiency at the system level (focusing on GPU memory IO efficiency), enabling training LLMs with longer context windows with the same computational budget.
 - Some researchers are proposing novel architectures to tackle this efficiency challenge, such as RWKV and RetNet.

Fine-Tuning in Deep Learning

Fine-tuning is a technique in deep learning that falls under the umbrella of transfer learning. It involves making slight adjustments to a pre-trained model's internal parameters to adapt it to a new, related task without starting the training process from scratch. Let's explore the significance, steps, and common strategies associated with fine-tuning.

Significance of Fine-Tuning

- **Efficiency:** Compared to training a model from scratch, fine-tuning significantly reduces time and resource requirements by building upon a pre-existing model's learned features.

- **Improved Performance:** Leveraging a pre-trained model's knowledge and representations enhances performance on task-specific datasets, leading to better results.
- **Data Efficiency:** Fine-tuning enables effective training even with limited labeled data, maximizing the utility of available resources.

Steps in Fine-Tuning

1. **Select a Pre-Trained Model:** Choose a pre-trained model that aligns with the nature of your task, ensuring compatibility with the required features.
2. **Adjust Model Architecture:** Modify the model's architecture, typically focusing on the top layers, to suit the specific requirements of your task. This may involve altering output neuron counts or layer configurations.
3. **Freeze or Unfreeze Layers:** Decide whether to freeze certain layers to retain their learned representations or unfreeze them to adapt to new data during fine-tuning, depending on the task complexity and dataset size.
4. **Training:** Train the modified model on the task-specific dataset, using a smaller learning rate compared to initial pre-training to prevent drastic changes to learned representations while allowing adaptation to new data.
5. **Fine-Tuning Strategies:** Experiment with hyperparameters, loss functions, and training strategies to optimize fine-tuning results iteratively.

Common Use Cases for Fine-Tuning

Fine-tuning is particularly effective in scenarios where:

- **Style Customization:** Tailoring models to reflect specific brand tones or qualitative aspects.
- **Specialization:** Adapting general models to domain-specific tasks, such as legal or medical contexts.
- **Few-Shot Learning:** Achieving strong performance with minimal labeled examples for specific classification tasks.

- **Addressing Edge Cases:** Ensuring model robustness by handling rare or unlikely scenarios appropriately.
- **Incorporating Proprietary Data:** Integrating company-specific datasets to enhance model relevance and performance.

Fine-Tuning Strategies

Full Fine-Tuning

- Update the entire neural network similarly to the pre-training process, demanding significant computational resources.

Partial Fine-Tuning

- Selectively update critical parameters while freezing others to reduce computational demands and stabilize the model.

Additive Fine-Tuning

- Introduce new parameters or layers to the model and train only these additional components while preserving the original pre-trained weights.

Transfer Learning in Deep Learning

Transfer learning is a fundamental technique in deep learning that involves leveraging knowledge gained from solving one task to improve performance on a related but different task. It's a valuable strategy when there's a scarcity of labeled data for the target task or when training resources are limited. Let's delve into the significance, process, and various aspects of transfer learning, akin to fine-tuning.

Significance of Transfer Learning

- **Efficiency:** Transfer learning accelerates model training by utilizing pre-existing knowledge, reducing the need for extensive retraining.
- **Improved Performance:** By transferring learned features and representations from a source task to a target task, models can achieve better performance on new datasets.

- **Data Efficiency:** Transfer learning maximizes the utility of available labeled data by repurposing it for related tasks, making it particularly beneficial in scenarios with limited data availability.

Process of Transfer Learning

1. **Select Source Task and Model:** Choose a source task and a pre-trained model that has been trained on a relevant dataset. The model should possess generalized features applicable to the target task.
2. **Feature Extraction or Fine-Tuning:** Decide between feature extraction and fine-tuning based on the similarity between the source and target tasks.
 - **Feature Extraction:** Extract high-level features learned by the pre-trained model and use them as input to a new model trained for the target task.
 - **Fine-Tuning:** Make minor adjustments to the pre-trained model's internal parameters to adapt it to the target task, while leveraging its learned representations.
3. **Adapt Model Architecture:** Modify the model's architecture as needed to align with the requirements of the target task, such as adjusting output neuron counts or layer configurations.
4. **Transfer Learning Strategies:** Experiment with different transfer learning strategies, including frozen feature extraction, fine-tuning specific layers, or applying domain adaptation techniques, to optimize model performance.

Common Use Cases for Transfer Learning

Transfer learning finds application in various scenarios, including:

- **Domain Adaptation:** Adapting models trained on one domain to perform well in a related but different domain.
- **Multi-Task Learning:** Training models to perform multiple tasks simultaneously by sharing knowledge across tasks.
- **Zero-Shot Learning:** Learning to perform tasks without any labeled examples by leveraging knowledge from related tasks.

- **Semi-Supervised Learning:** Combining labeled and unlabeled data to improve model performance, especially in scenarios with limited labeled data.

Transfer Learning Strategies

Frozen Feature Extraction

- Extract high-level features learned by a pre-trained model and use them as fixed inputs to a new model for the target task, preventing further updates to the pre-trained weights.

Fine-Tuning

- Adjust the internal parameters of a pre-trained model to adapt it to the target task while leveraging its learned representations, allowing for minor updates to the pre-trained weights.

Domain Adaptation

- Incorporate techniques such as adversarial training or domain-specific loss functions to align the learned representations of a model with the target domain.

Differences between Transfer Learning & Fine Tuning

Aspects	Transfer Learning	Full Fine-tuning
Definition	In transfer learning, only a small subset of the model's parameters or a few task-specific layers are trained while keeping the majority of the pre-trained model's parameters frozen. One of the	In full fine-tuning, all the parameters of a pre-trained model are updated during the training process. This means the weights of every layer in the model are adjusted based on the new training data. This

	popular transfer learning strategies is PEFT.	approach can lead to significant changes in the model's behavior and performance, tailored to the specific task or dataset at hand.
Objective	The goal is to adapt the pre-trained model to a new task with minimal changes to its parameters. This approach seeks to find an optimal balance between retaining the general knowledge acquired during pre-training and making enough task-specific adjustments to perform well on the new task.	The goal is to comprehensively adapt the entire pre-trained model to a new task or dataset. The aim is to maximize performance on that specific task.
Model Architecture	Utilizes existing architecture; Freezes most of the layers. Only a small set of parameters are fine-tuned.	With full fine-tuning, every parameter of the LLM gets updated.
Training Process	May involve only training a new top layer while keeping other layers fixed. In some cases, the number of newly trained parameters is just 1–2% of the original LLM weights.	Involves adjusting specific layers and parameters for the new task.
Data Requirement	Smaller dataset with fewer examples	Requires task-specific large data set for fine-tuning.
Computational Complexity	This approach is generally more resource-efficient, as only a small portion of the model is being updated. It requires less memory and processing power and often leads to faster training times. This makes it more accessible for situations with limited computational resources or for quick experimentation.	Since it involves updating all the parameters, full fine-tuning is typically more computationally intensive and time-consuming. It requires more memory and processing power, as well as potentially longer training times, especially for large models. Full fine-tuning often requires a large amount of GPU RAM.
Storage requirements	Reduced storage requirements	Increased storage requirements model

Model performance	Performance can be similar, but often a bit lower than full fine-tuning	Typically results in higher performance
Inference hosting Requirements	What is needed is to host original LLM and additional model weights for inference	Each fine-tuned model must be hosted

RAG vs Finetuning

Purpose:

- Fine-tuning LLMs: Fine-tuning involves taking a pre-trained LLM (such as GPT-3 or BERT) and adapting it to specific tasks. It is a versatile approach used for various NLP tasks, including text classification, language translation, sentiment analysis, and more. Fine-tuned LLMs are typically used when the task can be accomplished with just the model itself and doesn't require external information retrieval.
- RAG: RAG models are designed for tasks that involve both retrieval and generation of text. They combine a retrieval mechanism, which fetches relevant information from a large database, with a generation mechanism

that generates human-like text based on the retrieved information. RAG models are often used for question-answering, document summarization, and other tasks where accessing local information is crucial.

Architecture:

- **Fine-tuning LLMs:** Fine-tuning LLMs usually start with a pre-trained model (like GPT-3) and fine-tune it by training on task-specific data. The architecture remains largely unchanged, with adjustments made to the model's parameters to optimize performance for the specific task.
- **RAG:** RAG models have a hybrid architecture that combines a transformer-based LLM (like GPT) with an external memory module that allows for efficient retrieval from a knowledge source, such as a database or a set of documents.

Training Data:

- **Fine-tuning LLMs:** Fine-tuning LLMs rely on task-specific training data, often consisting of labeled examples that match the target task, but they don't explicitly involve retrieval mechanisms.
- **RAG:** RAG models are trained to handle both retrieval and generation, which typically involves a combination of supervised data (for generation) and data that demonstrates how to retrieve and use external information effectively

Use Cases:

- **Fine-tuning LLMs:** Fine-tuned LLMs are suitable for a wide range of NLP tasks, including text classification, sentiment analysis, text generation, and more, where the task primarily involves understanding and generating text based on the input.
- **RAG:** RAG models excel in scenarios where the task requires access to external knowledge, such as open-domain question-answering, document

summarization, or chatbots that can provide information from a knowledge base.

Data Collection and Preparation

- **Data Collection:** Collect data from various sources including internet, social media, sensors, and optimally data warehouses like Google BigQuery or Google Cloud Storage. When collecting data to train generative AI models, ensure that the dataset represents a diverse range of perspectives, backgrounds, and sources. Be cautious about relying solely on data from a single source or domain, as it may introduce bias.
- **Data Creation:** Optionally generate/use synthetic data to augment existing datasets to achieve better performance metrics on trained models.
- **Data Curation:** Clean and organize data through filtering, transformation, integration and labeling of data for supervised learning. Programmatic labeling may be an indispensable part of how we can accelerate training and fine-tuning.
- **Feature Storage:** Store the gold standard of features that have been painstakingly curated and prepped. We do this primarily for AI governance and consistency across the enterprise and the may project that will be attached to it; we will want to store and manage features used in training and inference ; often distinguishing between static and dynamic or train time and run time feature store capabilities. Looking to provide answers to questions from development teams around: "Where did we get the data from? Who touched the data? Did they bias the data during curation/labeling?" For example, you can use the [Vertex AI Feature Store](#) to support this task.
- **Data Bias Check.** This is an essential step in the machine learning lifecycle, especially for generative AI models, to ensure that the training data and the

generated content is fair, unbiased, and representative of the underlying data distribution. Bias in generative AI models can lead to the creation of outputs that reinforce stereotypes, discriminate against certain groups, or misrepresent reality. We can address data bias in predictive and generative AI, right after collecting and curating the data. By incorporating *data bias checks* into the machine learning lifecycle for generative AI, we can work towards creating models that generate fair, unbiased, and representative content, ultimately ensuring that the AI systems we develop align with ethical guidelines and societal values.

Data Collection and Preparation

Model Training & Experimentation

- Experimentation: This is name of the game in machine learning.
- Pre-trained model training: Train a model using a pre-trained model as a starting point
- Fine-tune models: Update an existing model with new data or for a new task
- Model registry: This is a means and mechanism for Model Governance and Model Risk Management. To use the Model Registry we would Store and manage models for versioning, reuse, and auditing, marking the meta-data of the model to see where the data came from (Data Provenance, Lineage), how the model was trained, etc.

Distributed Training

- Data parallelism: Split data across multiple devices and training them simultaneously to speed up training
- Model parallelism: Split the model across multiple devices and training them simultaneously to fit larger models into memory

- Federated learning: Train models collaboratively across multiple devices or organizations while preserving privacy and security.

Model Adaptation

Transfer learning

Adapt a pre-trained model to a new domain or task

Fine-tune models

Update an existing model with new data or for a new task

AI Safety Mitigation

Ensure the model is safe, ethical, and trustworthy by mitigating misuse, improving robustness, and ensuring privacy and security.

Distillation

Train or create a much smaller model with similar performance and downstream task capabilities by using a student/teacher few-shot learning approach to train a much smaller model, or, prune a model to arrive at a smaller footprint deployable model.

There are several methods to distill large language models (LLMs) to arrive at significantly smaller models with comparable performance. Let's explore some common approaches .

1. **Pruning:** This involves removing unimportant weights or neurons from the LLM while preserving its accuracy. Pruning can be done using various techniques such as magnitude-based pruning, sensitivity-based pruning, and weight-rewinding.
2. **Quantization:** This involves reducing the precision of the LLM's weights and activations. For example, weights that were initially represented as 32-bit floating-point numbers can be quantized to 8-bit integers. This reduces the model's memory footprint without significantly impacting its performance.
3. **Knowledge Distillation:** This involves training a smaller model to mimic the output of the LLM on a set of training examples. The idea is to transfer the knowledge from the larger model to the smaller one. This approach can be

used to train models that are orders of magnitude smaller than the original model with only a small reduction in accuracy.

4. **Low-rank Factorization:** This involves decomposing the weight matrices of the LLM into low-rank matrices. This reduces the number of parameters in the model while preserving its accuracy.
5. **Compact Embeddings:** This involves reducing the dimensionality of the input and output embeddings of the LLM. This reduces the number of parameters in the model and speeds up inference.
6. **Architectural Changes:** This involves changing the architecture of the LLM to make it more efficient. For example, the transformer architecture can be modified to reduce its memory footprint or the number of attention heads can be reduced.
7. **Parameter Sharing:** This involves sharing parameters between different parts of the model. For example, the same weight matrix can be used for multiple layers in a neural network, reducing the number of parameters and improving the model's efficiency.

Model Deployment

Let's discuss how to serve the model using [Google Cloud's Vertex AI](#) endpoint. Serving a model on a Vertex AI endpoint involves packaging the model into a container image, uploading the image to a container registry, creating a model resource and endpoint in the Vertex AI console, deploying the model to the endpoint, testing the endpoint, and monitoring its performance.

1. **Prepare the Model.** First, you need to prepare the trained model for deployment. This typically involves packaging the model and any required dependencies into a container image, such as a Docker image.
2. **Upload the Container Image.** Next, you need to upload the container image to a container registry, such as Google Container Registry. This will allow you to easily deploy the image to a Vertex AI endpoint.

3. **Create a Model Resource.** In the Vertex AI console, create a model resource that will represent your deployed model. This involves specifying the model's name, description, and any other metadata that may be required.
4. **Create an Endpoint.** Create a Vertex AI endpoint to serve your model. This involves specifying the endpoint's name, description, and the region where it will be deployed.
5. **Deploy the Model.** Once the endpoint is created, you can deploy your model to the endpoint. This involves specifying the container image to use and any required settings, such as the number of replicas to deploy.

Model Maintenance

- Data and model drift. Monitor and managing changes in data distribution and model performance over time
- Model versioning. Manage different versions of the model for traceability, reproducibility, and comparison
- Prompt-tuning, Prompt design: Optimize prompts used in language generation models to improve quality, coherence, and relevance.

Performance optimization

- Hardware acceleration. Optimize the model for specific hardware accelerators to improve inference speed and efficiency
- Quantization. Reduce the model size and computation requirements by converting it to lower precision data types
- Pruning. Reduce the model size and computation requirements by removing unimportant weights and neurons.