

EXP 1 Comprehensive Report on the Fundamentals of Generative AI and Large Language Models (LLMs)

AIM:

Comprehensive Report on the Fundamentals of Generative AI and Large Language Models (LLMs)

ALGORITHM:

Introduction

Foundations of Generative AI

Understanding Large Language Models (LLMs)

Training and Dataset Considerations

Applications of Generative AI and LLMs

Challenges and Limitations

What is generative AI?

Generative AI, sometimes called gen AI, is artificial intelligence (AI) that can create original content such as text, images, video, audio or software code in response to a user's prompt or request.

Generative AI relies on sophisticated machine learning models called deep learning models algorithms that simulate the learning and decision-making processes of the human brain. These models work by identifying and encoding the patterns and relationships in huge amounts of data, and then using that information to understand users' natural language requests or questions and respond with relevant new content.

AI has been a hot technology topic for the past decade, but generative AI, and specifically the arrival of ChatGPT in 2022, has thrust AI into worldwide headlines and launched an unprecedented surge of AI innovation and adoption. Generative AI offers enormous productivity benefits for individuals and organizations, and while it also presents very real challenges and risks, businesses are forging ahead, exploring how the technology can improve their internal workflows and enrich their products and services. According to research by the management consulting firm McKinsey, one third of organizations are already using generative AI regularly in at least one business function.¹ Industry analyst Gartner projects more than 80% of organizations will have deployed generative AI applications or used

generative AI application programming interfaces (APIs) by 2026.²

How generative AI works:

For the most part, generative AI operates in three phases:

Training, to create a foundation model that can serve as the basis of multiple gen AI applications.

Tuning, to tailor the foundation model to a specific gen AI application.

Generation, evaluation and retuning, to assess the gen AI application's output and continually improve its quality and accuracy.

Training

Generative AI begins with a foundation model, a deep learning model that serves as the basis for multiple different types of generative AI applications. The most common foundation models today are large language models (LLMs), created for text generation applications, but there are also foundation models for image generation, video generation, and sound and music generation as well as multimodal foundation models that can support several kinds content generation.

To create a foundation model, practitioners train a deep learning algorithm on huge volumes of raw, unstructured, unlabeled data e.g., terabytes of data culled from the internet or some other huge data source. During training, the algorithm performs and evaluates millions of 'fill in the blank' exercises, trying to predict the next element in a sequence e.g., the next word in a sentence, the next element in an image, the next command in a line of code and continually adjusting itself to minimize the difference between its predictions and the actual data (or 'correct' result).

The result of this training is a neural network of parameters, encoded representations of the entities, patterns and relationships in the data, that can generate content autonomously in response to inputs, or prompts.

This training process is compute-intensive, time-consuming and expensive: it requires thousands of clustered graphics processing units (GPUs) and weeks of processing, all of which costs millions of dollars. Open-source foundation model projects, such as Meta's Llama-2, enable gen AI developers to avoid this step and its costs.

Tuning

Metaphorically speaking, a foundation model is a generalist: It knows a lot about a lot of types of content, but often can't generate specific types of output with desired accuracy or fidelity. For that, the model must be tuned to a specific content generation task. This can be done in a variety of ways.

Fine tuning

Fine tuning involves feeding the model labeled data specific to the content generation application questions or prompts the application is likely to receive, and corresponding correct answers in the desired format. For example, if a development team is trying to create a customer service chatbot, it would create hundreds or thousands of documents containing labeled customers service questions and correct answers, and then feed those documents to the model.

Fine-tuning is labor-intensive. Developers often outsource the task to companies with large data-labeling workforces.

Reinforcement learning with human feedback (RLHF)

In RLHF, human users respond to generated content with evaluations the model can use to update the model for greater accuracy or relevance. Often, RLHF involves people 'scoring' different outputs in response to the same prompt. But it can be as simple as having people type or talk back to a chatbot or virtual assistant, correcting its output.

Generation, evaluation, more tuning

Developers and users continually assess the outputs of their generative AI apps, and further tune the model even as often as once a week for greater accuracy or relevance. (In contrast, the foundation model itself is updated much less frequently, perhaps every year or 18 months.)

Another option for improving a gen AI app's performance is retrieval augmented generation (RAG). RAG is a framework for extending the foundation model to use relevant sources outside of the training data, to supplement and refine the parameters or representations in the original model. RAG can ensure that a generative AI app always has access to the most current information. As a bonus, the additional sources accessed via RAG are transparent to users in a way that the knowledge in the original foundation model is not.

Generative AI is poised to affect multiple roles within various sectors:

 Business and Legal <ul style="list-style-type: none"> • Purchasing agents • Compensation specialists • Management analysts • Market research analysts • Marketing specialists • Lawyers and paralegals 	 Finance <ul style="list-style-type: none"> • Insurance underwriters • Budget analysts • Accountants and auditors • Personal financial advisors • Credit professionals • Financial analysts • Tax preparers 	 Social Sciences <ul style="list-style-type: none"> • Geographers • Epidemiologists • Survey researchers • Political scientists • Sociologists • Economists 	 Writing and Editing <ul style="list-style-type: none"> • Writers and authors • Reporters and correspondents • Technical writers • Interpreters and translators • Editors
 STEM <ul style="list-style-type: none"> • Programmers and software developers • Web developers • Some types of engineers • Data scientists • Physicists • Medical scientists • Operations research analysts 	 Sales <ul style="list-style-type: none"> • Insurance sales agents • Advertising sales agents • Travel agents • Securities, commodities and financial sellers • Telemarketers 	 Office and Administrative Support <ul style="list-style-type: none"> • Procurement clerks • Credit authorizers, checkers and clerks • Cargo and freight agents • Statistical assistants • Loan interviewers and clerks • Billing and posting clerks 	 Other <ul style="list-style-type: none"> • Postsecondary teachers • Public relations specialists • Interior designers

Generative AI model architectures and how they have evolved

Truly generative AI models deep learning models that can autonomously create content on demand have evolved over the last dozen years or so. The milestone model architectures during that period include

Variational autoencoders (VAEs), which drove breakthroughs in image recognition, natural language processing and anomaly detection.

Generative adversarial networks (GANs) and diffusion models, which improved the accuracy of previous applications and enabled some of the first AI solutions for photo-realistic image generation.

Transformers, the deep learning model architecture behind the foremost foundation models and generative AI solutions today.

Variational autoencoders (VAEs)

An autoencoder is a deep learning model comprising two connected neural networks: One that encodes (or compresses) a huge amount of unstructured, unlabeled training data into parameters, and another that decodes those parameters to reconstruct the content. Technically, autoencoders can generate new content, but they're more useful for compressing data for storage or transfer, and decompressing it for use, than they are for high-quality content generation.

Introduced in 2013, variational autoencoders (VAEs) can encode data like an autoencoder, but decode multiple new variations of the content. By training a VAE to generate variations toward a particular goal, it can 'zero in' on more accurate, higher-fidelity content over time. Early VAE applications included anomaly detection (e.g., medical image analysis) and natural language generation.

Generative adversarial networks (GANs)

GANs, introduced in 2014, also comprise two neural networks: A generator, which generates new content, and a discriminator, which evaluates the accuracy and quality the generated data. These adversarial algorithms encourages the model to generate increasingly high-quality outputs.

GANs are commonly used for image and video generation, but can generate high-quality, realistic content across various domains. They've proven particularly successful at tasks as style transfer (altering the style of an image from, say, a photo to a pencil sketch) and data augmentation (creating new, synthetic data to increase the size and diversity of a training data set).

Diffusion models

Also introduced in 2014, diffusion models work by first adding noise to the training data until it's random and unrecognizable, and then training the algorithm to iteratively diffuse the noise to reveal a desired output.

Diffusion models take more time to train than VAEs or GANs, but ultimately offer finer-grained control over output, particularly for high-quality image generation tool. DALL-E, Open AI's image-generation tool, is driven by a diffusion model.

Transformers

First documented in a 2017 paper published by Ashish Vaswani and others, transformers evolve the encoder-decoder paradigm to enable a big step forward in the way foundation models are trained, and

in the quality and range of content they can produce. These models are at the core of most of today's headline-making generative AI tools, including ChatGPT and GPT-4, Copilot, BERT, Bard, and Midjourney to name a few.

Transformers use a concept called attention, determining and focusing on what's most important about data within a sequence to;

process entire sequences of data e.g., sentences instead of individual words simultaneously;

capture the context of the data within the sequence;

encode the training data into embeddings (also called hyperparameters) that represent the data and its context.

In addition to enabling faster training, transformers excel at natural language processing (NLP) and natural language understanding (NLU), and can generate longer sequences of data e.g., not just answers to questions, but poems, articles or papers with greater accuracy and higher quality than other deep generative AI models. Transformer models can also be trained or tuned to use tools e.g., a spreadsheet application, HTML, a drawing program to output content in a particular format.

What generative AI can create:

Generative AI can create many types of content across many different domains.

Text-

Generative models, especially those based on transformers, can generate coherent, contextually relevant text, everything from instructions and documentation to brochures, emails, web site copy, blogs, articles, reports, papers, and even creative writing. They can also perform repetitive or tedious writing tasks (e.g., such as drafting summaries of documents or meta descriptions of web pages), freeing writers' time for more creative, higher-value work.

Images and video-

Image generation such as DALL-E, Midjourney and Stable Diffusion can create realistic images or original

art, and can perform style transfer, image-to-image translation and other image editing or image enhancement tasks. Emerging gen AI video tools can create animations from text prompts, and can apply special effects to existing video more quickly and cost-effectively than other methods.

Sound, speech and music-

Generative models can synthesize natural-sounding speech and audio content for voice-enabled AI chatbots and digital assistants, audiobook narration and other applications. The same technology can generate original music that mimics the structure and sound of professional compositions.

Software code-

Gen AI can generate original code, autocomplete code snippets, translate between programming languages and summarize code functionality. It enables developers to quickly prototype, refactor, and debug applications while offering a natural language interface for coding tasks.

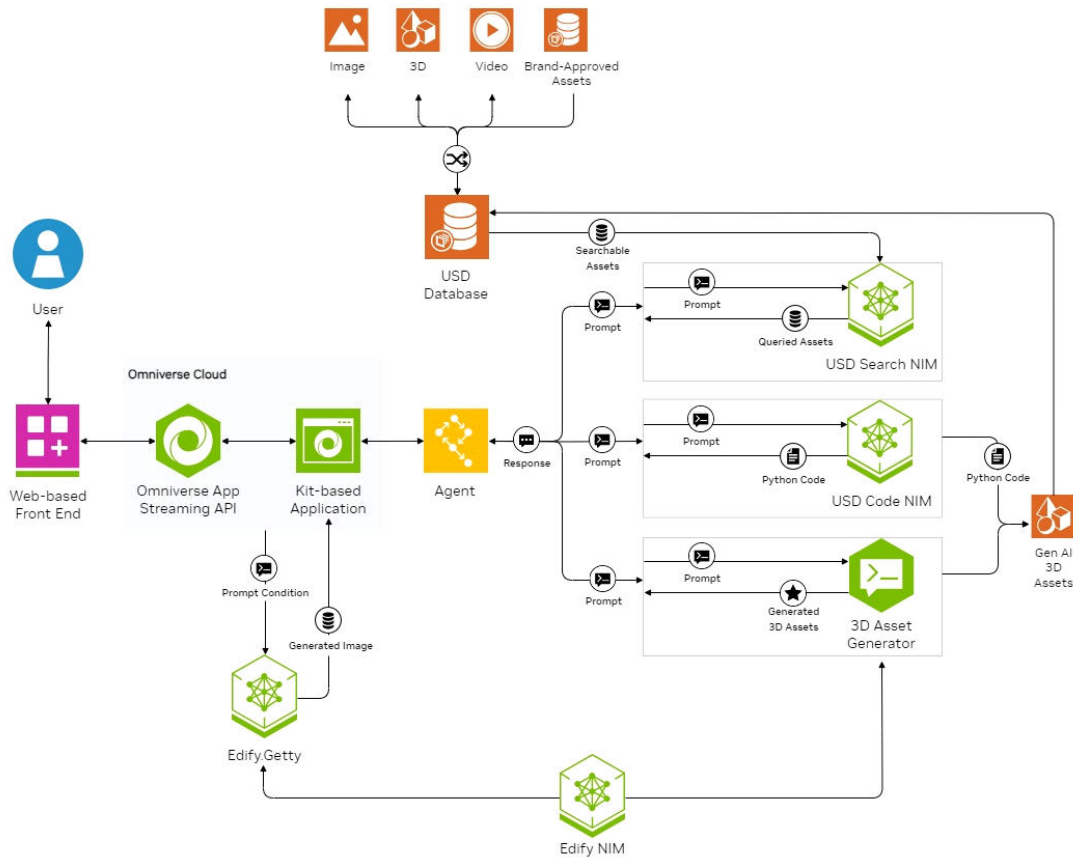
Design and art-

Generative AI models can generate unique works of art and design, or assist in graphic design. Applications include dynamic generation of environments, characters or avatars, and special effects for virtual simulations and video games.

Simulations and synthetic data-

Generative AI models can be trained to generate synthetic data, or synthetic structures based on real or synthetic data. For example, generative AI is applied in drug discovery to generate molecular structures with desired properties, aiding in the design of new pharmaceutical compounds.

Architecture of generative ai



Benefits of generative AI

The obvious, overarching benefit of generative AI is greater efficiency. Because it can generate content and answers on demand, gen AI has the potential to accelerate or automate labor-intensive tasks, cut costs, and free employees time for higher-value work.

But generative AI offers several other benefits for individuals and organizations.

Enhanced creativity

Gen AI tools can inspire creativity through automated brainstorming, generating multiple novel versions of content. These variations can also serve as starting points or references that help writers, artists, designers and other creators plow through creative blocks.

Improved (and faster) decision-making

Generative AI excels at analyzing large datasets, identifying patterns and extracting meaningful insights, then generating hypotheses and recommendations based on those insights to support executives,

analysts, researchers and other professionals in making smarter, data-driven decisions.

Dynamic personalization

In applications like recommendation systems and content creation, generative AI can analyze user preferences and history and generate personalized content in real time, leading to a more tailored and engaging user experience.

Constant availability

Generative AI operates continuously without fatigue, providing around-the-clock availability for tasks like customer support chatbots and automated responses.

Use cases for generative AI:

The following are just a handful of gen AI use cases for enterprises. As the technology develops and organizations embed these tools into their workflows, we can expect to see many more.

Customer experience

Marketing organizations can save time and amp up their content production by using gen AI tools to draft copy for blogs, web pages, collateral, emails and more. But generative AI solutions can also produce highly personalized marketing copy and visuals in real time based on when, where and to whom the ad is delivered. And it will power next-generation chatbots and virtual agents that can give personalized responses and even initiate actions on behalf of customer, a significant advancement compared to the previous generation of conversational AI models trained on more limited data for very specific tasks.

Software development and application modernization

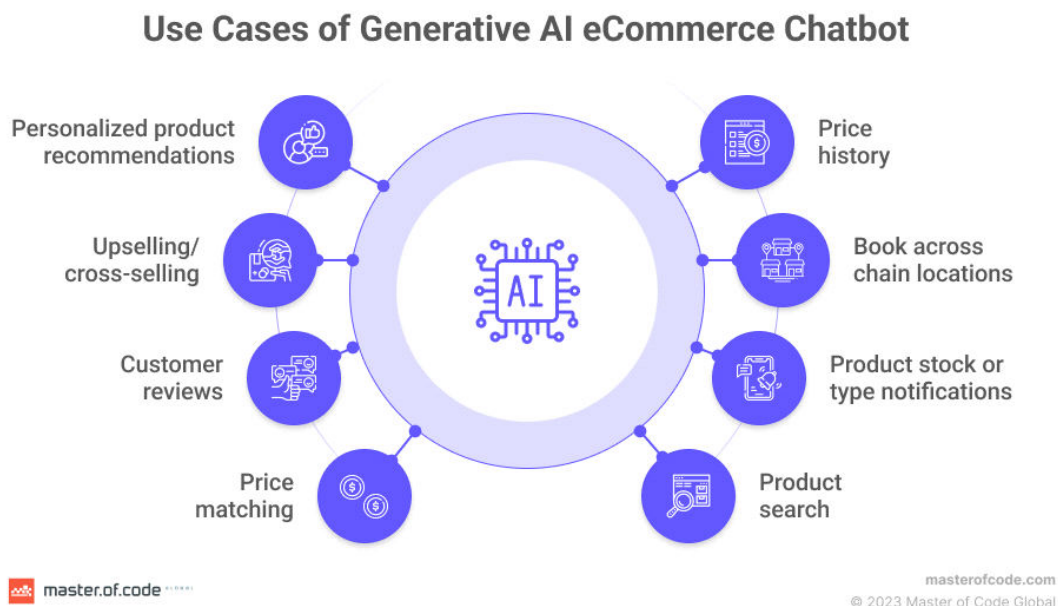
Code generation tools can automate and accelerate the process of writing new code. Code generation also has the potential to dramatically accelerate application modernization by automating much of the repetitive coding required to modernize legacy applications for hybrid cloud environments.

Digital labor

Generative AI can quickly draw up or revise contracts, invoices, bills and other digital or physical 'paperwork' so that employees who use or manage it can focus on higher level tasks. This can accelerate workflows in virtually every enterprise area including human resources, legal, procurement and finance.

Science, engineering and research

Generative AI models can help scientists and engineers propose novel solutions to complex problems. In healthcare, for example, generative models can be applied to synthesize medical images for training and testing medical imaging systems.



Challenges, limitations and risks

Generative AI has made remarkable strides in a relatively short period of time, but still presents significant challenges and risks to developers, users and the public at large. Below are some of the most serious issues, and how they're being addressed.

'Hallucinations' and other inaccurate outputs

An AI hallucination is a generative AI output that is nonsensical or altogether inaccurate but, all too often, seems entirely plausible. The classic example is when a lawyer used a gen AI tool for research in preparation for a high-profile case and the tool 'produced' several example cases, complete with quotes and attributions, that were entirely fictional.

Some practitioners view hallucinations as an unavoidable consequence of balancing a model's accuracy and its creative capabilities. But developers may implement preventative measures, called guardrails,

that restrict the model to relevant or trusted data sources. Continual evaluation and tuning can also help reduce hallucinations and inaccuracies.

Inconsistent outputs

Due to the variational or probabilistic nature of gen AI models, the same inputs can result in slightly or significantly different outputs. This can be undesirable in certain applications, such as customer service chatbots, where consistent outputs are expected or desired. Through prompt engineering iteratively refining or compounding prompts, users can arrive at prompts that consistently deliver the results they want from their generative AI applications.

Bias

Generative models may learn societal biases present in the training data or in the labeled data, external data sources, or human evaluators used to tune the model and generate biased, unfair or offensive content as a result. To prevent biased outputs from their models, developers must ensure diverse training data, establish guidelines for preventing bias during training and tuning, and continually evaluate model outputs for bias as well as accuracy.

Lack of explainability and metrics

Many generative AI models are 'black box' models, meaning it can be challenging or impossible to understand their decision-making processes; even the engineers or data scientists who create the underlying algorithm can understand or explain what exactly is happening inside it and how it arrives at a specific result. Explainable AI practices and techniques can help practitioners and users understand and trust the processes and outputs of generative models.

Assessing and comparing the quality of generated content can also be challenging. Traditional evaluation metrics may not capture the nuanced aspects of creativity, coherence, or relevance. Developing robust and reliable evaluation methods for generative AI remains an active area of research.

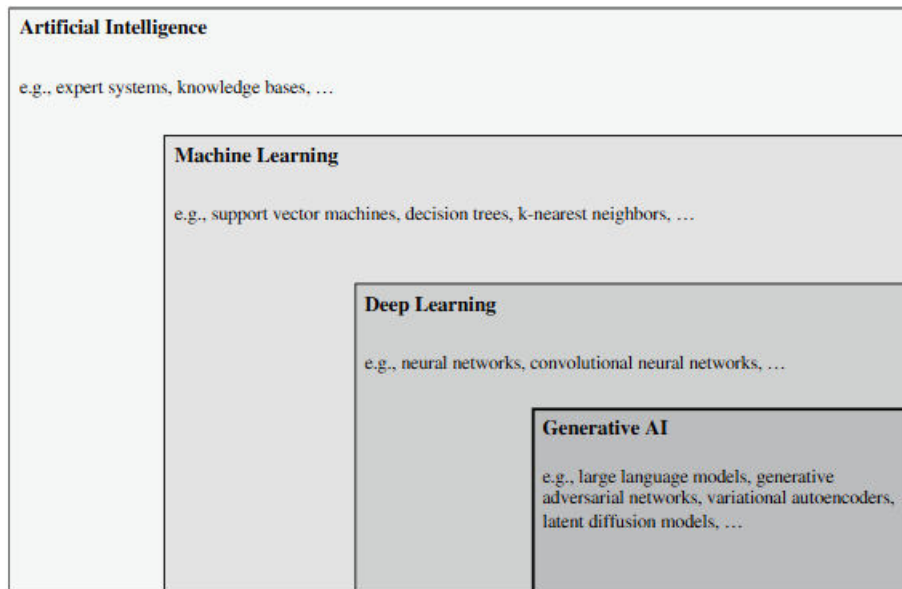
Threats to security, privacy and intellectual property

Generative AI models can be exploited to generate convincing phishing emails, fake identities or other malicious content that can fool users into taking actions that compromise security and data privacy. Developers and users need to be careful that data put into the model (during tuning, or as part of a prompt) doesn't expose their own intellectual property (IP) or any information protected as IP by other organizations. And they need to monitor outputs for new content that exposes their own IP or violates others' IP protections.

Deepfakes

Deepfakes are AI-generated or AI-manipulated images, video or audio created to convince people that they're seeing, watching or hearing someone do or say something they never did or said. They are among the most chilling examples of how the power of generative AI can be applied with malicious intent.

Most people are familiar with deepfakes created to damage reputations or spread misinformation. More recently, cybercriminals have deployed deepfakes as part of cyberattacks (e.g., fake voices in voice phishing scams) or financial fraud schemes.



Conceptual framework of generative AI

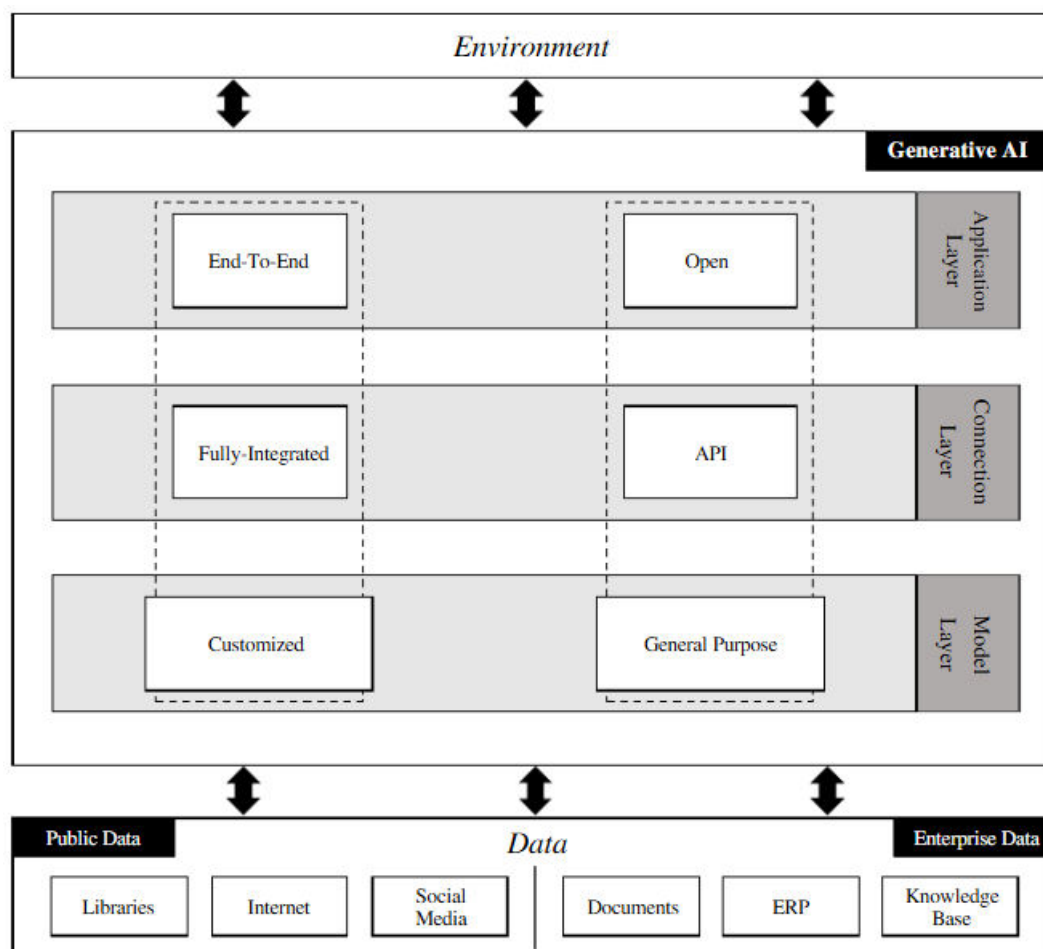


Table 1 Overview of core deep generative models

Deep generative model

Description

Generative adversarial network (GAN)

Description

Generative adversarial networks consist of two competing neural networks: a generator and a discriminator (Goodfellow et al., 2020). The generator creates realistic data samples, while the discriminator distinguishes between real and generated samples (Pan et al., 2019). Both neural networks are trained together until the discriminator is not able to differentiate both samples (Janiesch et al., 2021). This adversarial competition results in the generator improving its data generation capabilities over time, eventually producing high-quality, realistic outputs. Hence, GANs find various applications, for instance, in image generation and manipulation, object detection and segmentation, and natural language processing (Aggarwal et al., 2021; Gui et al., 2023)

Deep generative model

Variational autoencoder (VAE)

Description

Variational autoencoders employ a neural network to learn encoding compressed input data into a lower-dimensional latent space and then decode the data by reconstructing the original data from the latent space representation (Kingma et al., 2014). By optimizing a variational lower bound on the data likelihood in a probabilistic approach, VAEs can generate new samples that resemble the original data distribution. Typical use cases for VAEs can be seen in the synthetic generation and reconstruction of data such as images, in anomaly detection, and recommendation systems (Wei & Mahmood, 2021)

Deep generative model

Transformer

Description

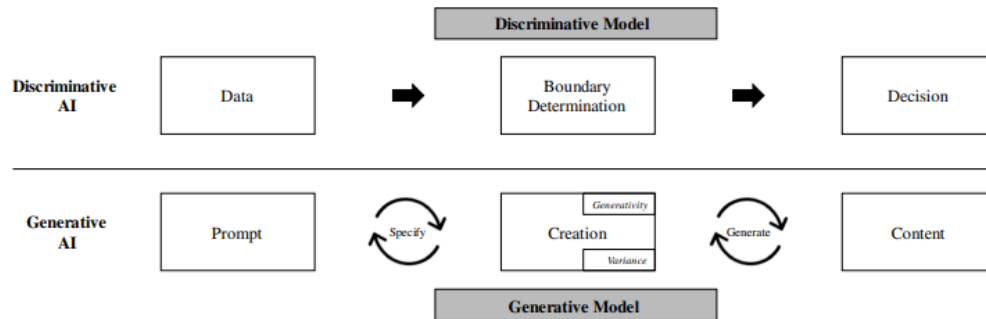
Transformer models have become the basis for many state-of-the-art natural language processing tasks and succeeding models. They are a specific type of neural network architecture that employ self-attention mechanisms to capture long-range dependencies in the data, making them well-suited for large-scale language modeling tasks (Vaswani et al., 2017). Generative pre-trained transformers (GPT) build on the transformer architecture and were trained with large datasets of unlabeled data (Brown et al., 2020). Due to their large size (i.e., a very large number of trainable parameters), GPT trained on text data are often referred to as large language models (LLMs) (Schramowski et al., 2022). The goal of LLMs is to generate novel, coherent, contextually relevant human-like text by predicting which token is most likely to occur after the prior tokens in a sentence (Brown et al., 2020; H. Li, 2022). Hence, LLMs can serve as the foundation for conversational AI tools like ChatGPT (Teubner et al., 2023). Besides conversing, the large amount of information stored in LLMs can be used for text generation, writing, or even programming, e.g., to support scholars (Cooper, 2023; Lund et al., 2023)

Latent diffusion model (LDM)

Description

Latent diffusion models are transformer based and build on the concepts of denoising score matching and contrastive divergence to learn a stochastic data generation process (Rombach et al., 2022). In LDMs, the generation process starts with a simple initial distribution, such as Gaussian noise. Then, the data gets gradually refined through a series of noise-reduction steps following a predefined diffusion process through a latent space (Ho et al., 2020). The key advantage of LDMs is their ability to learn complex data distributions without requiring adversarial training (as in GANs) or optimizing variational lower bounds (as in VAEs). They also feature improved stability over other DGMs during training to be less prone to issues like mode collapse (Kodali et al., 2017; Rombach et al., 2022), making them well-suited for high-quality and detailed outputs, such as high-resolution image synthesis (Ho et al., 2020)

Procedural differences of generative AI and discriminative AI



Overview of different output modalities for generative AI applications

Modality--Text

Description:

X-to-text applications are centered around text generation and natural language processing. The goal is to generate human-like written text that fits the user's input prompt by providing a meaningful answer within the context. For instance, chatbots like OpenAI's ChatGPT imitate textual conversations with the user and can be guided to output text artifacts as desired (OpenAI, 2023). Furthermore, text-producing applications can be leveraged for content creation (e.g., copywriting or specific writing in e-commerce contexts) (Bakpayev et al., 2022; Brand et al., 2023). Moreover, text generation can support processes in sales or support by providing the ability to produce customized texts tailored toward the requests (Mondal et al., 2023). Systems integrating GAI models with further knowledge bases (e.g., enterprise data and Internet access) extend the available information beyond the model's initial training dataset.

Modality--Image

Description:

X-to-image applications generate images based on the user's prompting. Relying on GANs or diffusion models as DGMs, synthetic images are created that find use cases in marketing, design and fashion, or creative fields in the form of new visual art (Haase et al., 2023; Mayahi & Vidrih, 2022; Zhang et al., 2023a). For instance, Stable Diffusion is an open-sourced x-to-image model that enables the generation of images in multiple GAI applications (Rombach et al., 2022). Moreover, generated synthetic images can act as training data for further ML models to train classifiers (e.g., medical images to detect diseases (Ali et al., 2023)). Besides a text-to-image creation process, image editing capabilities are possible, e.g., via image-to-image systems that manipulate and extend images according to the user's prompting (Openlaender, 2022).

Modality--Video

Description:

X-to-video applications deal with the creation of synthetic videos, i.e., dynamic motion images. New video clips are generated by describing the content of the desired video footage (text-to-video) or applying the style and composition via text or image prompt to a source video (video-to-video) (Esser et al., 2023). These prospects allow the fast and convenient creation and editing of videos via natural language and other modalities (Zhan et al., 2021). Thus, not only videographers benefit from x-to-video applications but also people without filming and editing skills are enabled to creatively express themselves due to an accessible creation process (Anantrasirichai & Bull, 2022). Besides recreational and

entertainment purposes, x-to-video GAI models find application in sales and marketing (e.g., product marketing videos), onboarding and education (e.g., virtual avatars in training videos), or in customer support (e.g., how-to videos) (Leiker et al., 2023; Mayahi & Vidrih, 2022). As an exemplary application, Synthesia is a video creation platform specialized in generating professional videos with virtual avatars and synthetic voiceovers (Synthesia, 2023)

Modality--Code

Description:

In the realm of software development, x-to-code GAI applications offer transformative potential in how developers work and code by providing x-to-text capabilities specific to programming languages. Models like CodeBERT (Feng et al., 2020) or GraphCodeBERT (Guo et al., 2021) were trained on programming code to generate source code from natural language or modeling languages for new software programs. Several x-to-text models also offer coding capabilities because general-purpose LLMs are trained with increasingly large datasets that contain code (e.g., Stability.ai, 2023). Programmers using applications such as GitHub Copilot are supported by automatically written chunks of code, ideas converted into actionable scripts, auto-completion functions, generated unit tests, duplicate code detection, and bug fixing (Sun et al., 2022). These automation potentials allow developers to focus on higher-level tasks and problem solving, enhancing their productivity and the final product's overall quality, reducing time-to-market, supporting rapid prototyping, and promoting continuous innovation for the product and business

Modality--Audio

Description:

X-to-audio applications focus on audio content generation and comprise, for instance, the generation of speech with synthetically generated human-like voices (Borsos et al., 2022; Wang et al., 2023). Especially text-to-speech and speech-to-speech models are being heavily researched and can be used to power various applications, ranging from digital assistants and customer services to audiobook and training narration and accessibility tools (Moussawi et al., 2021; Qiu & Benbasat, 2005). GAI models like Microsoft's VALL-E (Wang et al., 2023) offer a more personalized and engaging user experience by enabling realistic voice modeling. Moreover, x-to-sound models find application in music creation. By specifying genres or melodies via prompts, unique pieces of music can be generated that respect the original intent (Agostinelli et al., 2023). GAI models such as MusicLM (Agostinelli et al., 2023) help musicians in their creative process, offering inspiration and aiding the composition of complex pieces. Businesses in the music industry can leverage high-fidelity music generation to create customized soundtracks for marketing, movies, or video games, significantly reducing the cost and time associated with traditional music production (Anantrasirichai & Bull, 2022; Weng & Chen, 2020)

Modality--Other

Description:

The applications of GAI extend beyond the stated modality types and domains, impacting multiple other, specific areas. For instance, x-to-molecules models like AlphaFold (Jumper et al., 2021) and OpenBioML (Murphy & Thomas, 2023) generate viable protein structures and design new molecules by generating valid, novel molecular structures, supporting drug discovery and bioengineering researchers (Walters & Murcko, 2020). 3D modeling is also impacted by GAI applications such as DreamFusion (Poole et al., 2023), Nvidia GET3D (Gao et al., 2022), and Point-E (Nichol et al., 2022), which generate realistic and complex 3D models that facilitate a range of applications from product design and architecture to virtual reality and game development

A brief history of generative AI:

The term “generative AI” exploded into the public consciousness in the 2020s, but gen AI has been part of our lives for decades, and today’s generative AI technology draws on machine learning breakthroughs from as far back as the early 20th century. A non-exhaustive representative history of generative AI might include some of the following dates

1964: MIT computer scientist Joseph Weizenbaum develops ELIZA, a text-based natural language processing application. Essentially the first chatbot (called a ‘chatterbot’ at the time), ELIZA used pattern-matching scripts to respond to typed natural language inputs with empathetic text responses.

1999: Nvidia ships GeForce, the first graphical processing unit. Originally developed to deliver smooth motion graphics for video games, GPUs had become the defacto platform for developing AI models and mining cryptocurrencies.

2004: Google autocomplete first appears, generating potential next words or phrases as users enter their search terms. The relatively modern example of generative AI is based on a Markov Chain, a mathematical model developed in 1906.

2013: The first variational autoencoders (VAEs) appear.

2014: The first generative adversarial networks (GANs) and diffusion models appear.

2017: Ashish Vaswani, a team at Google Brain, and a group from the University of Toronto publish “Attention is All You Need,” a paper documenting the principles of transformer models, widely acknowledged as enabling the most powerful foundation models and generative AI tools being developed today.

2019-2020: OpenAI rolls out its GPT (Generative Pretrained Transformer) large language models, GPT-2 and GPT-3.

2022: OpenAI introduces ChatGPT, a front-end to GPT-3 that generates complex, coherent and contextual sentences and long-form content in response to end-user prompts.

What are LLMs?

Large language models (LLMs) are a category of foundation models trained on immense amounts of data making them capable of understanding and generating natural language and other types of content to perform a wide range of tasks.

LLMs have become a household name thanks to the role they have played in bringing generative AI to the forefront of the public interest, as well as the point on which organizations are focusing to adopt artificial intelligence across numerous business functions and use cases.

Outside of the enterprise context, it may seem like LLMs have arrived out of the blue along with new developments in generative AI. However, many companies, including IBM, have spent years implementing LLMs at different levels to enhance their natural language understanding (NLU) and natural language processing (NLP) capabilities. This has occurred alongside advances in machine learning, machine learning models, algorithms, neural networks and the transformer models that provide the architecture for these AI systems.

LLMs are a class of foundation models, which are trained on enormous amounts of data to provide the foundational capabilities needed to drive multiple use cases and applications, as well as resolve a multitude of tasks. This is in stark contrast to the idea of building and training domain specific models for each of these use cases individually, which is prohibitive under many criteria (most importantly cost and infrastructure), stifles synergies and can even lead to inferior performance.

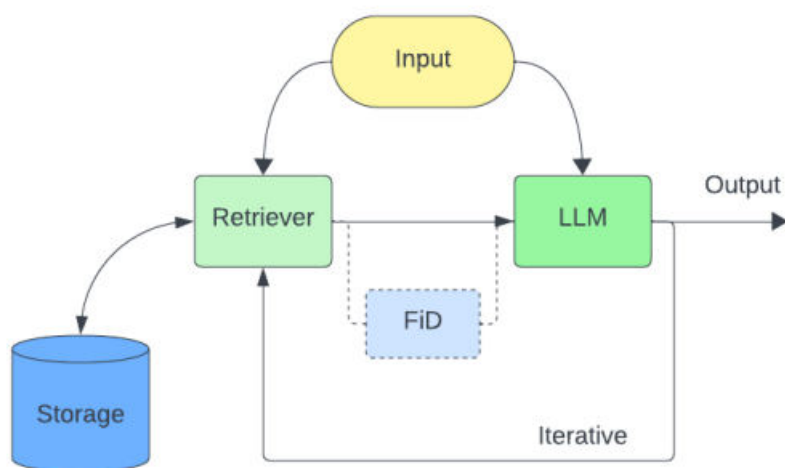
LLMs represent a significant breakthrough in NLP and artificial intelligence, and are easily accessible to the public through interfaces like Open AI's Chat GPT-3 and GPT-4, which have garnered the support of Microsoft. Other examples include Meta's Llama models and Google's bidirectional encoder representations from transformers (BERT/RobERTa) and PaLM models. IBM has also recently launched its Granite model series on watsonx.ai, which has become the generative AI backbone for other IBM products like watsonx Assistant and watsonx Orchestrate.

In a nutshell, LLMs are designed to understand and generate text like a human, in addition to other forms of content, based on the vast amount of data used to train them. They have the ability to infer from context, generate coherent and contextually relevant responses, translate to languages other than English, summarize text, answer questions (general conversation and FAQs) and even assist in creative writing or code generation tasks.

They are able to do this thanks to billions of parameters that enable them to capture intricate patterns in language and perform a wide array of language-related tasks. LLMs are revolutionizing applications in various fields, from chatbots and virtual assistants to content generation, research assistance and language translation.

As they continue to evolve and improve, LLMs are poised to reshape the way we interact with technology and access information, making them a pivotal part of the modern digital landscape.

A flow diagram of Retrieval Augmented LLMs. The retriever extracts a similar context to the input and forwards it to the LLM either in simple language or encoded through Fusion-in-Decoder (FiD). Depending on the task, retrieval and generation may repeat multiple times.



How large language models work:

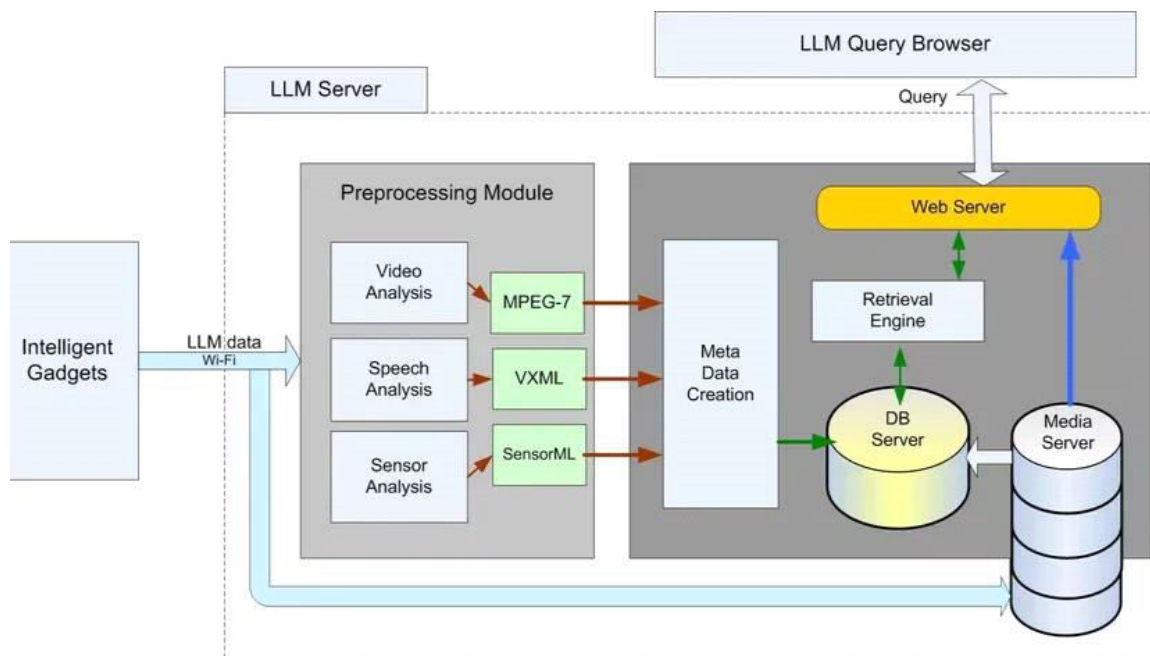
LLMs operate by leveraging deep learning techniques and vast amounts of textual data. These models are typically based on a transformer architecture, like the generative pre-trained transformer, which excels at handling sequential data like text input. LLMs consist of multiple layers of neural networks, each with parameters that can be fine-tuned during training, which are enhanced further by a numerous layer known as the attention mechanism, which dials in on specific parts of data sets.

During the training process, these models learn to predict the next word in a sentence based on the context provided by the preceding words. The model does this through attributing a probability score to the recurrence of words that have been tokenized, broken down into smaller sequences of characters. These tokens are then transformed into embeddings, which are numeric representations of this context.

To ensure accuracy, this process involves training the LLM on a massive corpora of text (in the billions of pages), allowing it to learn grammar, semantics and conceptual relationships through zero-shot and self-supervised learning. Once trained on this training data, LLMs can generate text by autonomously predicting the next word based on the input they receive, and drawing on the patterns and knowledge they've acquired. The result is coherent and contextually relevant language generation that can be harnessed for a wide range of NLU and content generation tasks.

Model performance can also be increased through prompt engineering, prompt-tuning, fine-tuning and other tactics like reinforcement learning with human feedback (RLHF) to remove the biases, hateful speech and factually incorrect answers known as “hallucinations” that are often unwanted byproducts of training on so much unstructured data. This is one of the most important aspects of ensuring enterprise-grade LLMs are ready for use and do not expose organizations to unwanted liability, or cause damage to their reputation.

Architecture of IIm:



LLM use cases :

LLMs are redefining an increasing number of business processes and have proven their versatility across a myriad of use cases and tasks in various industries. They augment conversational AI in chatbots and virtual assistants (like IBM watsonx Assistant and Google’s BARD) to enhance the interactions that underpin excellence in customer care, providing context-aware responses that mimic interactions with

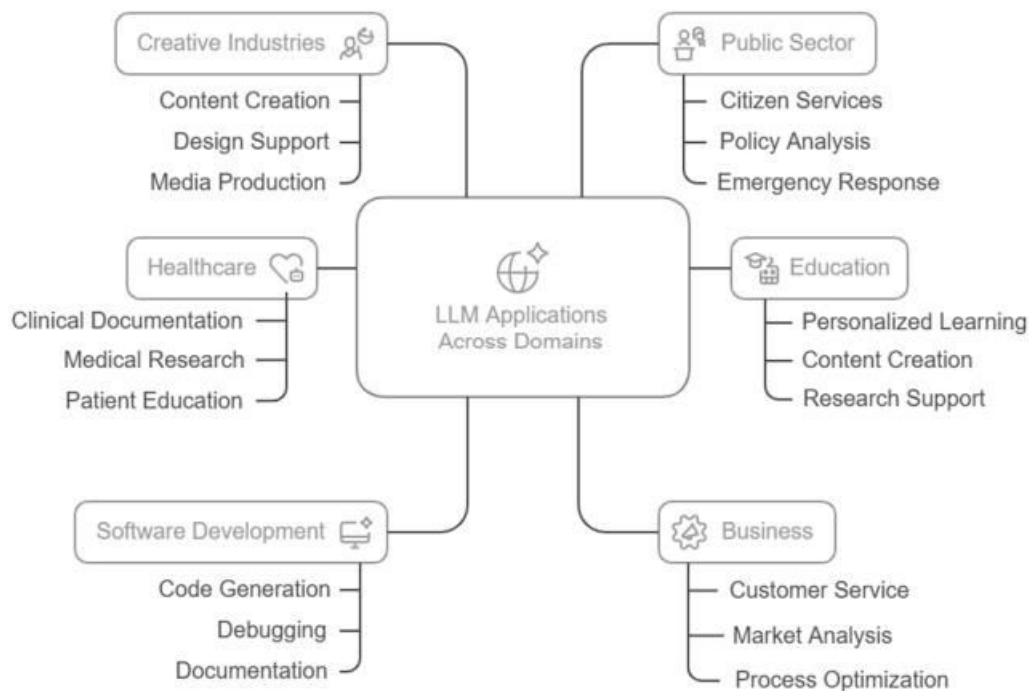
human agents.

LLMs also excel in content generation, automating content creation for blog articles, marketing or sales materials and other writing tasks. In research and academia, they aid in summarizing and extracting information from vast datasets, accelerating knowledge discovery. LLMs also play a vital role in language translation, breaking down language barriers by providing accurate and contextually relevant translations. They can even be used to write code, or “translate” between programming languages.

Moreover, they contribute to accessibility by assisting individuals with disabilities, including text-to-speech applications and generating content in accessible formats. From healthcare to finance, LLMs are transforming industries by streamlining processes, improving customer experiences and enabling more efficient and data-driven decision making.

Most excitingly, all of these capabilities are easy to access, in some cases literally an API integration away.

LLM applications in different domains:



Here is a list of some of the most important areas where LLMs benefit organizations:

Text generation: language generation abilities, such as writing emails, blog posts or other mid-to-long

form content in response to prompts that can be refined and polished. An excellent example is retrieval-augmented generation (RAG).

Content summarization: summarize long articles, news stories, research reports, corporate documentation and even customer history into thorough texts tailored in length to the output format.

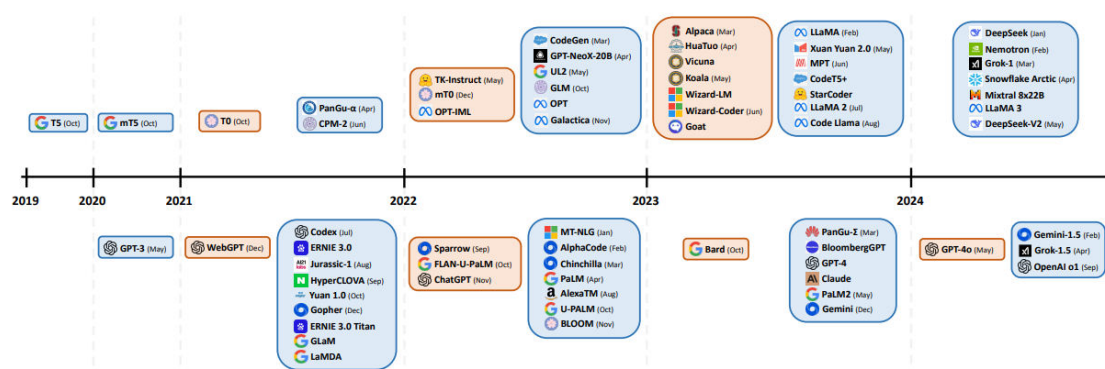
AI assistants: chatbots that answer customer queries, perform backend tasks and provide detailed information in natural language as a part of an integrated, self-serve customer care solution.

Code generation: assists developers in building applications, finding errors in code and uncovering security issues in multiple programming languages, even “translating” between them.

Sentiment analysis: analyze text to determine the customer’s tone in order understand customer feedback at scale and aid in brand reputation management.

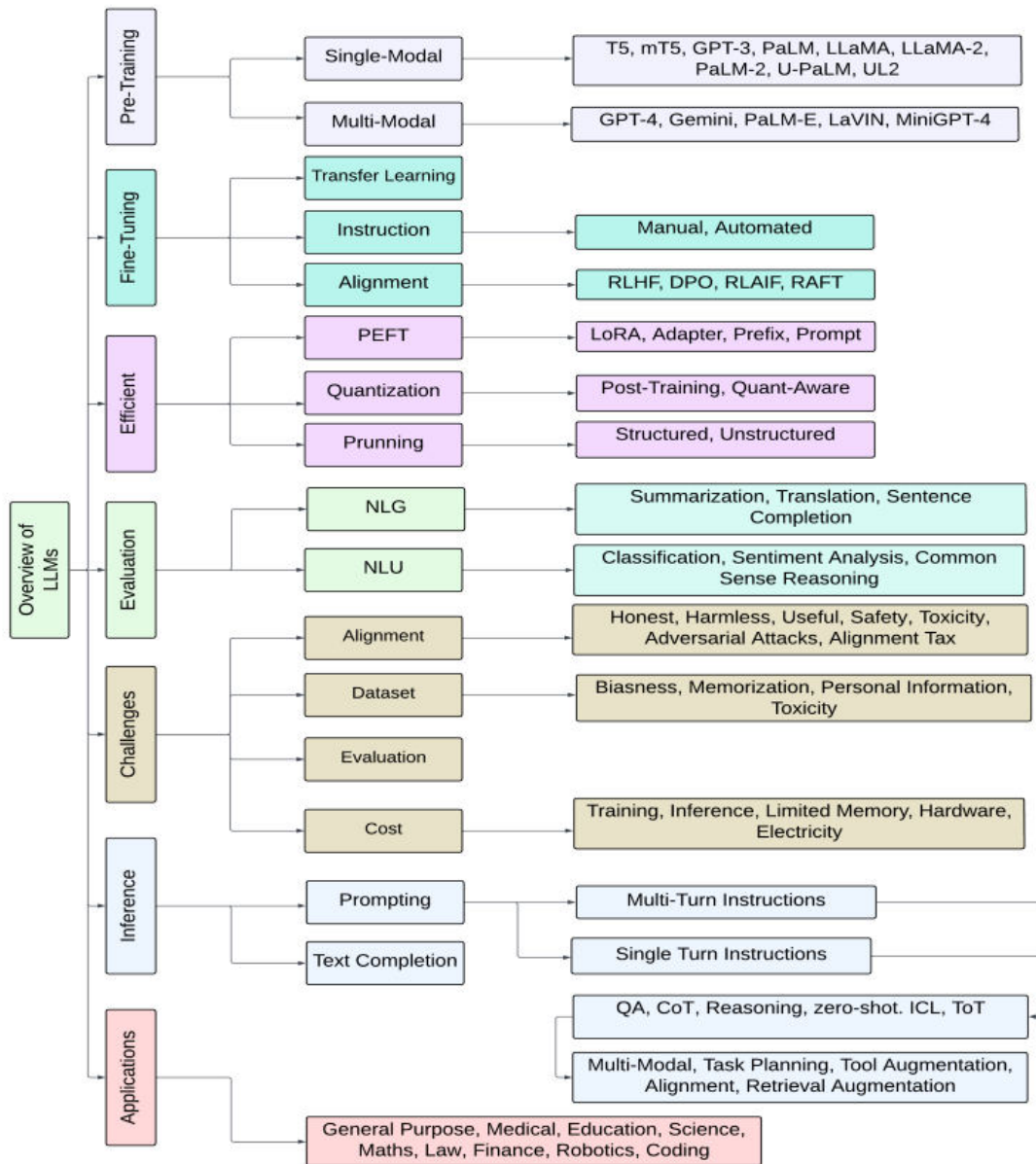
Language translation: provides wider coverage to organizations across languages and geographies with fluent translations and multilingual capabilities.

LLMs stand to impact every industry, from finance to insurance, human resources to healthcare and beyond, by automating customer self-service, accelerating response times on an increasing number of tasks as well as providing greater accuracy, enhanced routing and intelligent context gathering.

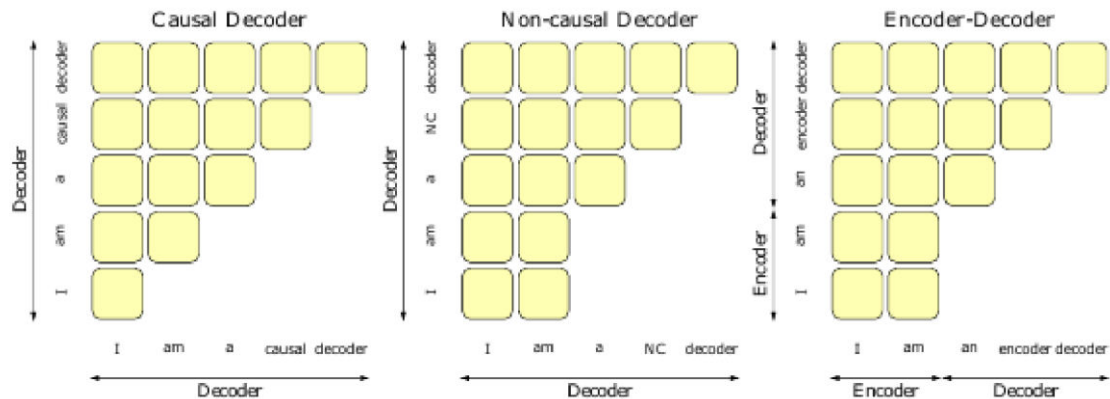


Chronological display of LLM releases: blue cards represent ‘pre-trained’ models, while orange cards correspond to ‘instruction-tuned’ models. Models on the upper half signify open-source availability, whereas those on the bottom are closed-source. The chart illustrates the increasing trend towards instruction-tuned and open-source models, highlighting the evolving landscape and trends in natural language processing research.

A broader overview of LLMs, dividing LLMs into seven branches: 1. Pre-Training 2. Fine-Tuning 3. Efficient 4. Inference 5. Evaluation 6. Applications 7. Challenges



An example of attention patterns in language models:



Key insights and findings from the study of instruction-tuned Large Language Models:

Models--T0

- Multi-task prompting enables zero-shot generalization and outperforms baselines
- Even a single prompt per dataset task is enough to improve performance

Models--WebGPT

- To aid the model in effectively filtering and utilizing relevant information, human labelers play a crucial role in answering questions regarding the usefulness of the retrieved documents
- Interacting a fine-tuned language model with a text-based web-browsing environment can improve end-to-end retrieval and synthesis via imitation learning and reinforcement learning
- Generating answers with references can make labelers easily judge the factual accuracy of answers

Models--Tk-INSTRUCT

- Instruction tuning leads to a stronger generalization of unseen tasks
- More tasks improve generalization whereas only increasing task instances does not help
- Supervised trained models are better than generalized models
- Models pre-trained with instructions and examples perform well for different types of inputs

Models--mT0 and BLOOMZ

- Instruction tuning enables zero-shot generalization to tasks never seen before
- Multi-lingual training leads to even better zero-shot generalization for both English and nonEnglish
- Training on machine-translated prompts improves performance for held-out tasks with non-English prompts

- English only fine-tuning on multilingual pre-trained language model is enough to generalize to other pre-trained language tasks

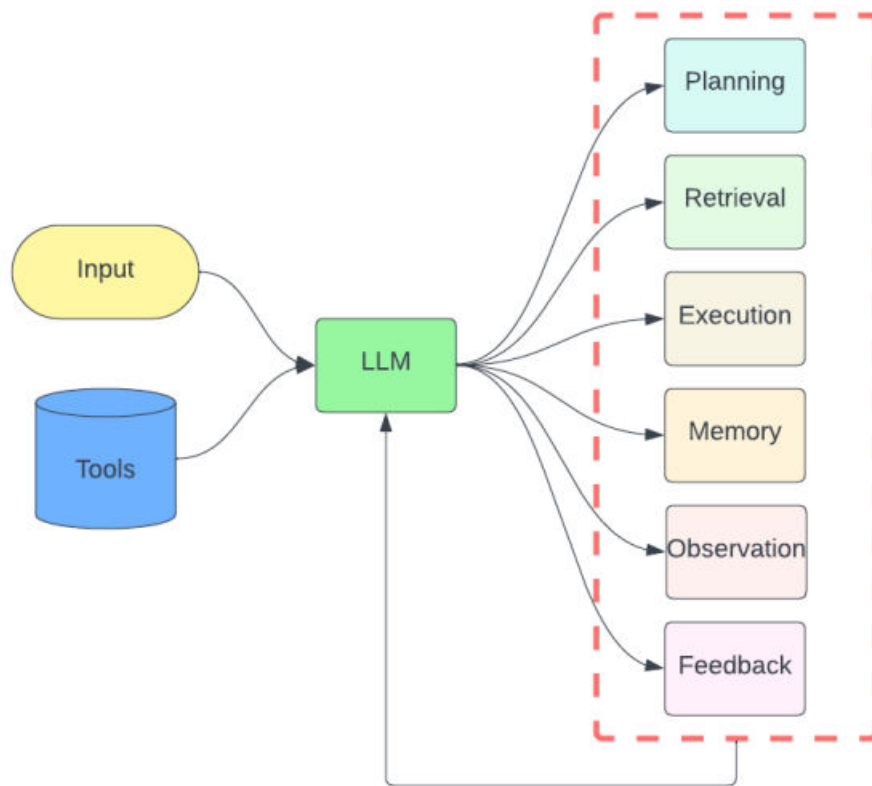
Models--OPT-IML

- Creating a batch with multiple task examples is important for better performance
- Only example proportional sampling is not enough, training datasets should also be proportional for better generalization/performance
- Fully held-out and partially supervised tasks performance improves by scaling tasks or categories whereas fully supervised tasks have no effect
- Including small amounts i.e. 5% of pretraining data during fine-tuning is effective
- Only 1% reasoning data improves the performance, adding more deteriorates performance
- Adding dialogue data makes the performance worse

Models--Sparrow

- Labelers' judgment and well-defined alignment rules help the model generate better responses
- Good dialogue goals can be broken down into detailed natural language rules for the agent and the raters
- The combination of reinforcement learning (RL) with reranking yields optimal performance in terms of preference win rates and resilience against adversarial probing

A basic flow diagram of tool augmented LLMs.:



Challenges of LLMs

LLMs pose a number of concerns.

Like broader generative AI technologies, LLMs pose significant threats to certain industries, such as finance, journalism, and customer support.

In education and academia, LLMs can enable individuals to cheat on assignments and papers. According to Nature, there have been numerous papers published in journals with the phrase “regenerate response” — indicating the text was copied from an LLM like ChatGPT.

Data bias is another major challenge, as these LLMs often replicate and amplify biases from their training data. A study published by Apple's Machine Learning Research on four different LLMs highlighted how these models are prone to stereotyping professions based on gender.

A particularly pressing legal challenge involves copyright infringement issues. LLMs typically require vast amounts of data for training, which they often get by scraping content from various sources, including copyrighted materials without explicit permission. Recently the New York Times and several other US newspapers sued OpenAI and Microsoft for copyright infringement, highlighting the complex ethical and

legal landscape that surrounds LLMs.

LLMs also can generate and spread misleading or false information, posing risks to information integrity and public trust.

Generative AI vs. LLMs Recap:

Generative Capabilities--

Both Generative AI and LLMs are capable of generating new content.

Generative AI can produce a variety of content types like images, videos, and text.

LLMs are a subset of generative AI that specializes in generating coherent, contextually relevant text.

Core Technologies--

Generative AI uses technologies like Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs). These models learn to create new outputs by mimicking the distribution of input data.

LLMs (Large Language Models) use transformer models. Transformers use self-attention to weigh the relevance of all parts of the text to each other. This makes LLMs effective for tasks that require a deep understanding of language.

Data Usage--

Generative AI models require diverse and large datasets to effectively create new content.

LLMs specifically require large volumes of high-quality text data.

Application Areas--

Generative AI has broad applications across many industries, including creative fields, science, finance, and more.

LLMs excel in environments that demand high levels of text interaction, such as customer support systems and educational tools. Additionally, LLMs are used in industries like finance for tasks like fraud detection, where they analyze textual data to identify anomalies.

Ethical and Practical Challenges--

Both LLMs and Generative AI deal with data bias and copyright concerns due to their reliance on extensive datasets.

Generative AI poses unique challenges with the potential creation of deepfakes.

LLMs have been criticized for enabling academic dishonesty and potentially spreading misinformation due to their ability to generate convincing textual content.

Final Thoughts--

Generative AI encompasses a broad range of technologies, including Large Language Models (LLMs).

While generative AI as a whole pushes the boundaries of creative content production, LLMs specifically refine how we generate and interact with textual data.

The integration of these technologies into various sectors brings transformative potential but also poses significant ethical challenges and risks.

By understanding the specific roles and capabilities within the broader spectrum of generative AI, we can better navigate these technologies and use them effectively.

Conclusion:

It contributes to summarizing significant findings of LLMs in the existing literature and provides a detailed analysis of the design aspects, including architectures, datasets, and training pipelines. We identified crucial architectural components and training strategies employed by different LLMs. These aspects are presented as summaries and discussions throughout the article. Moreover, we have discussed the performance differences of LLMs in zero-shot and few-shot settings, explored the impact of fine-tuning, and compared supervised and generalized models and encoder vs. decoder vs. encoder-decoder architectures. A comprehensive review of multi-modal LLMs, retrieval augmented LLMs, LLMs-powered agents, efficient LLMs, datasets, evaluation, applications, and challenges is also provided. This article is anticipated to serve as a valuable resource for researchers, offering insights into the recent advancements in LLMs and providing fundamental concepts and details to develop better LLMs.