

What is Generative AI?

Generative AI refers to a class of artificial intelligence models capable of producing novel and realistic content, such as text, images, audio, and video, that resembles data they were trained on. Unlike traditional AI that might analyze or classify existing data, Generative AI *creates* new data. Its core function is to learn the underlying patterns and structures of input data to then generate new samples that adhere to those learned distributions.



Types of GenAI Models

Generative AI models are typically categorized by the type of data they are designed to generate:

T

Text Models:

- Generate human-like text, including articles, stories, code, summaries, and conversational responses.
- Examples: GPT series (GPT-3, GPT-4), BERT (used for generation tasks in some contexts), T5, LLaMA, Mistral.



Image Models:

- Create realistic or artistic images from text descriptions (text-to-image), generate variations of existing images, or fill in missing parts of images.
- Examples: DALL-E, Midjourney, Stable Diffusion, GANs (Generative Adversarial Networks).



Audio Models:

- Generate speech (text-to-speech), music, sound effects, or even synthesize new voices.
- Examples: WaveNet, VALL-E, MusicGen.



Video Models:

- Generate short video clips from text descriptions, animate still images, or create realistic motion.
- Examples: Sora, Make-A-Video, RunwayML Gen-2.

Difference between Generative and Discriminative Models

This is a fundamental distinction in machine learning:

Generative Models:

- **Goal:** To learn the *distribution* of the input data and often the joint probability distribution $P(X, Y)$ (where X is input and Y is output).
- **Functionality:** Can *generate* new data samples that resemble the training data. They model how data is created.
- **Tasks:** Data generation, anomaly detection, completing missing data.
- **Examples:** GANs, VAEs, Autoregressive Models (like GPT for text generation).

Discriminative Models:

- **Goal:** To learn the *conditional probability distribution* $P(Y|X)$ (i.e., given an input X , what is the probability of output Y). They focus on the decision boundary between classes.
- **Functionality:** Can *discriminate* or classify between different categories or predict an output given an input. They do not generate new data.
- **Tasks:** Classification, regression.
- **Examples:** Support Vector Machines (SVMs), Logistic Regression, standard Neural Networks for classification, Decision Trees.

Transformers and Attention Mechanism (Brief Overview)

Transformers are a groundbreaking neural network architecture that revolutionized sequence modeling, particularly in Natural Language Processing (NLP). The core innovation of the Transformer is the **Attention Mechanism**.

Transformers:

- Introduced in the "Attention Is All You Need" paper (2017).
- **Key Idea:** Rely entirely on attention mechanisms to draw global dependencies between input and output, rather than using recurrent or convolutional layers. This allows for parallel processing of input sequences, making them significantly faster and more scalable for long sequences.
- **Architecture:** Typically consists of an encoder (for understanding input) and a decoder (for generating output), though some models (like GPT) are decoder-only. Both contain multiple layers of self-attention and feed-forward networks.

Attention Mechanism:

- **Core Function:** Allows the model to weigh the importance of different parts of the input sequence when processing a specific element. Instead of processing tokens sequentially, attention allows the model to "look at" and "attend to" all other tokens in the sequence simultaneously to understand context.
- **Self-Attention:** A specific type of attention where the attention is computed within the same sequence. For example, when generating a word, the model can attend to all previous words in the input to understand their relationship and contribute to the current word's meaning.
- **Mechanism:** It works by computing query, key, and value vectors from the input embeddings. The dot product of queries and keys determines "attention scores," which are then used to weight the values, producing a context-aware representation for each element.

Popular GenAI Models

These models represent key milestones and widely used architectures in Generative AI:



GPT (Generative Pre-trained Transformer):

- Developed by OpenAI.
- **Architecture:** Decoder-only Transformer.
- **Key Feature:** Trained on vast amounts of text data to predict the next word in a sequence. This "pre-training" allows them to learn grammar, facts, reasoning abilities, and writing styles.
- **Use Cases:** Text generation, translation, summarization, question answering, chatbots.



BERT (Bidirectional Encoder Representations from Transformers):

- Developed by Google.
- **Architecture:** Encoder-only Transformer.
- **Key Feature:** Trained bidirectionally, meaning it considers the context from both the left and right sides of a word, enabling a deeper understanding of language.
- **Use Cases:** Primarily for understanding language (classification, question answering, sentiment analysis), but its representations can be used for downstream generation tasks.



T5 (Text-to-Text Transfer Transformer):

- Developed by Google.
- **Architecture:** Encoder-decoder Transformer.
- **Key Feature:** Frames *all* NLP tasks as a "text-to-text" problem, where the input is text and the output is also text. This unified approach simplifies model architecture and training.
- **Use Cases:** Translation, summarization, question answering, text generation.



LLaMA (Large Language Model Meta AI):

- Developed by Meta AI.
- **Key Feature:** Family of foundation large language models (LLMs) known for being highly performant and often more accessible for research and fine-tuning due to their smaller (relative to some GPT models) yet powerful architectures.
- **Use Cases:** Similar to GPT, for various text-based generation and understanding tasks.



Mistral:

- Developed by Mistral AI.
- **Key Feature:** Another family of highly performant and efficient LLMs, often lauded for their compact size and strong performance, especially for their parameter count. They incorporate innovative techniques like Grouped Query Attention (GQA) for efficiency.
- **Use Cases:** General text generation, coding, summarization.

Fine-tuning vs. Prompt Engineering

These are two primary methods for customizing the behavior of pre-trained Generative AI models:

Fine-tuning:

- **Definition:** The process of taking a pre-trained model and continuing to train it on a smaller, specific dataset relevant to a particular task or domain. This involves updating the model's weights.
- **When to Use:** When you need the model to learn new knowledge, adapt to a specific style, or perform a specialized task that the base model isn't ideally suited for (e.g., generating medical reports, legal documents).
- **Pros:** Can achieve highly accurate and customized results; the model truly learns and integrates new information.
- **Cons:** Requires a dataset for training; computationally more expensive and time-consuming than prompt engineering; requires machine learning expertise.

Prompt Engineering:

- **Definition:** The art and science of crafting effective inputs (prompts) to guide a pre-trained Generative AI model to produce desired outputs without changing the model's underlying weights.
- **When to Use:** When you want to leverage the model's existing knowledge and abilities for various tasks; for quick iteration and experimentation; when you don't have a specific dataset for fine-tuning.
- **Pros:** Fast and flexible; no training required; accessible to non-technical users.
- **Cons:** Can be limited by the model's pre-trained knowledge; results can be inconsistent if prompts aren't well-designed; sometimes requires trial-and-error.

Ethics and Bias in Generative AI

Generative AI, like any powerful technology, comes with significant ethical considerations and potential for bias:

1	<p>Bias:</p> <ul style="list-style-type: none">• Source: Generative models learn from the data they are trained on. If this data contains societal biases (e.g., gender stereotypes, racial prejudice, historical inaccuracies), the model will learn and perpetuate these biases in its generated content.• Manifestations: Discriminatory language, stereotypes in generated images/text, unfair recommendations, perpetuating misinformation.• Mitigation: Diverse and balanced training datasets, explicit bias detection and mitigation techniques (e.g., debiasing algorithms), post-generation filtering, continuous monitoring.
2	<p>Misinformation and Disinformation:</p> <ul style="list-style-type: none">• Issue: GenAI can generate highly convincing but fabricated text, images, and videos (deepfakes), making it difficult to distinguish truth from falsehood. This can be used for malicious purposes like propaganda, scams, or discrediting individuals.• Mitigation: Watermarking, provenance tracking, improved detection methods for synthetic media, media literacy education.
3	<p>Copyright and Intellectual Property:</p> <ul style="list-style-type: none">• Issue: The use of copyrighted material in training data raises questions about fair use and compensation for creators. Also, who owns the copyright of AI-generated content?• Mitigation: Developing legal frameworks, exploring licensing models, using curated or open-source datasets.
4	<p>Job Displacement:</p> <ul style="list-style-type: none">• Issue: As GenAI becomes more capable, it may automate tasks traditionally performed by humans, leading to job displacement in creative industries, content creation, and more.• Mitigation: Focusing on re-skilling and up-skilling programs, fostering human-AI collaboration.
5	<p>Malicious Use:</p> <ul style="list-style-type: none">• Issue: GenAI can be used for cyberattacks (e.g., generating sophisticated phishing emails), creating harmful content, or impersonation.• Mitigation: Robust safety filters, access controls, responsible deployment.

Evaluation Metrics

Evaluating the quality and performance of Generative AI models is challenging because "good" generation is subjective. However, several metrics are commonly used:

For Text Generation:

- **BLEU (Bilingual Evaluation Understudy):**
 - **Purpose:** Measures the similarity of generated text to one or more reference texts, typically used for machine translation.
 - **Mechanism:** Calculates n-gram precision (overlap of sequences of words) between the generated text and reference text. Higher scores indicate greater similarity.
- **ROUGE (Recall-Oriented Understudy for Gisting Evaluation):**
 - **Purpose:** Measures the overlap between a generated summary and a set of reference summaries, typically used for text summarization.
 - **Mechanism:** Focuses on recall (how many n-grams in the reference are present in the generated text). Different ROUGE variants (ROUGE-N, ROUGE-L, ROUGE-S) exist.
- **Human Evaluation:**
 - **Purpose:** The most reliable way to assess subjective qualities like coherence, relevance, creativity, and factual accuracy.
 - **Mechanism:** Human annotators rate generated outputs based on predefined criteria.

General Considerations:

- **Diversity:** Does the model generate a wide variety of outputs, or does it tend to produce similar results repeatedly (mode collapse)?
- **Fidelity/Realism:** How realistic or faithful are the generated outputs to the training data distribution?
- **Coherence/Consistency:** Do the generated outputs make sense internally and logically?
- **Controllability:** How well can users control the characteristics of the generated output (e.g., style, content)?