# Learning Generative AI

## Introduction to Generative AI (GenAI)

Generative AI is a type of artificial intelligence that focuses on creating new, original content that resembles the data it was trained on. Instead of just analyzing or classifying existing data (like our Narrow AI examples), GenAI can produce things like text, images, audio, videos, and even code.

Think of it this way:

- A **discriminative AI** (a common type of Narrow AI) is like a detective who can identify a criminal based on evidence it has seen before. It learns to distinguish between different categories. For example, an image classifier that tells you if a picture contains a cat or a dog is a discriminative model.
- A **generative AI** is more like an artist who can create a brand new painting or a musician composing a new song based on the styles they've studied. It learns the underlying patterns in the data and then uses that knowledge to generate something novel.

Here are some examples of different GenAI models and what they can create:

- **Large Language Models (LLMs):** These are masters of text generation. They can write articles, poems, code, translate languages, answer your questions, and even have conversations (as you're experiencing now!). Examples include models like GPT-4, LaMDA, and others.
- **Image Generation Models:** These models can create realistic or artistic images from text descriptions or even modify existing images. Examples include DALL-E 2, Midjourney, and Stable Diffusion.
- **Audio Generation Models:** These can generate realistic speech, music, and sound effects.
- **Video Generation Models:** While still a rapidly evolving area, these models are starting to be able to create short video clips from text prompts.

While GenAI offers incredible potential, it also introduces new attack surfaces and amplifies existing security risks. For example, the very fact that LLMs can generate human-like text makes them potent tools for crafting sophisticated phishing attacks or spreading disinformation. Image and video generation can be exploited for creating convincing deepfakes for social engineering or even to bypass biometric authentication in the future.

Furthermore, the models themselves can be vulnerable. Think about:

- **Prompt Injection:** Maliciously crafted inputs can trick the model into bypassing safety measures or revealing sensitive information.
- **Data Poisoning:** If the training data is compromised, the model might learn to behave in undesirable or harmful ways.

- **Model Extraction:** Attackers might try to reverse-engineer the model to steal its intellectual property or gain insights into its workings to find vulnerabilities.
- **Supply Chain Risks:** The components and dependencies used to build and deploy AI models can also introduce security weaknesses.

At their heart, LLMs are a type of neural network, a computational system inspired by the human brain. What makes LLMs particularly powerful is their architecture, which is primarily based on something called the **Transformer** network.

Imagine you're trying to understand a sentence. You don't just look at each word in isolation; you pay attention to how the words relate to each other. The Transformer architecture allows LLMs to do something similar, but on a massive scale. It enables the model to understand the **context** of words in a sequence, no matter how far apart they are. This is a huge leap forward from earlier language models that struggled with long-range dependencies.

Here's a simplified way to think about the Transformer:

1. **Input:** The LLM takes a sequence of words (or parts of words) as input.
2. **Embeddings:** Each word is converted into a numerical representation called an "embedding." Think of this as a dense vector that captures the meaning of the word in a multi-dimensional space. Words with similar meanings will have embeddings that are closer together in this space.
3. **Attention Mechanisms:** This is the core innovation of the Transformer. The attention mechanism allows the model to weigh the importance of different words in the input sequence when processing a particular word. For example, if the model is trying to understand "The cat sat on the mat," the attention mechanism will help it realize that "sat" is closely related to both "cat" (the subject) and "mat" (the object).
4.
5. **Layers:** The Transformer network consists of multiple layers of these attention mechanisms and other processing units. These layers work together to progressively understand the relationships and patterns in the input text.
6. **Output:** Finally, the model generates an output, which could be the next word in a sequence, an answer to a question, a translation, or any other text-based task it has been trained for.

LLMs are typically trained in two main stages:

- **Pre-training:** This involves training the model on a massive dataset of text from the internet (think books, articles, websites, etc.). During pre-training, the model learns the basic structure of language, grammar, common knowledge, and relationships between words. A common pre-training task is "next word prediction" – the model learns to predict the next word in a sentence.
- **Fine-tuning:** After pre-training, the model can be further trained on a smaller, more specific dataset to perform particular tasks, like question answering, text summarization, or even generating code. This fine-tuning process helps the model specialize its general language understanding for specific applications.

Imagine you're at a crowded party, and someone across the room shouts your name. Even with all the background noise, your ears (and your brain!) can focus on that specific sound because it's important to you. The attention mechanism in an LLM works in a somewhat similar way.

When an LLM is processing a sentence, say "The dog chased the ball," and it's trying to understand the word "chased," the attention mechanism allows it to "pay more attention" to the words "dog" (the subject) and "ball" (the object) because they are highly relevant to the action of chasing. It assigns different "weights" or levels of importance to each word in the input sequence when processing every other word.

Think of it like this: for the word "chased," the model might give a high "attention score" to "dog" and "ball," a medium score to "the," and a lower score to other words if the sentence were longer. This helps the model understand the relationships between words and the overall meaning of the sentence.

Now, for a simplified visual of embeddings and the Transformer:

Imagine we have the sentence: "A happy cat sleeps."

**Embeddings:** Each word gets turned into a list of numbers (the embedding). Think of it like a secret code for each word that captures its meaning.

```
A:      [0.2, 0.5, -0.1, ...]
happy:  [0.8, 0.1, 0.3, ...]
cat:    [0.1, 0.9, -0.2, ...]
sleeps: [0.5, 0.3, 0.7, ...]
```

1. These numbers capture relationships. For example, the embedding for "cat" might be closer in numerical space to the embedding for "kitten" than to the embedding for "car."

**Transformer (Simplified):**

Imagine a series of processing boxes.

```
[Embedding(A)] --> [Attention] --> [Further Processing] -->
[Embedding(happy)] --> [Attention] --> [Further Processing] -->
[Embedding(cat)] --> [Attention] --> [Further Processing] -->
[Embedding(sleeps)] --> [Attention] --> [Further Processing] -->
                                                       ↓
                                                    [Output
(e.g., the next word)]
```

2. Inside the "Attention" box for the word "cat," the model looks at the embeddings of all the other words ("A," "happy," "sleeps") and assigns them scores based on how relevant they are to "cat." For instance, "happy" might get a higher score than "A" in understanding the context of "cat." These scores are then used to weigh the

information from the other words, helping the model understand the relationships within the sentence. This processed information then moves through further layers to make predictions or understand the input.

# Key Concepts Related to LLMs

## Tokens

LLMs don't process text as whole words. Instead, they break it down into smaller units called "tokens." A token can be a whole word, a part of a word, or even a single character. For example, the word "unbelievable" might be tokenized as "un", "believe", "able". The way text is tokenized can affect the model's performance and even its vulnerability to certain types of attacks. Think of tokens as the fundamental building blocks that the LLM understands and manipulates.

Why do LLMs use tokens instead of whole words? There are a few reasons:

- **Handling Rare Words:** If a model only worked with whole words, it would need to have a representation for every single word in the language, including rare or misspelled ones. This would lead to a massive vocabulary size. Tokenization allows the model to break down rare words into more common sub-word units. For example, "unbelievable" can be broken into "un", "believe", and "able", which are more frequent and have their own learned representations.
- **Flexibility:** Tokenization allows the model to handle words it hasn't seen before. If a new word appears, as long as its sub-word components exist in the model's vocabulary, it can still process it to some extent.
- **Efficiency:** Working with smaller units can sometimes be more computationally efficient.

The way text is tokenized can vary between different LLMs. Some models might use whole words for common words and sub-word units for rarer ones. Others might primarily use sub-word units based on frequency. Common tokenization methods include Byte-Pair Encoding (BPE) and WordPiece.

Here's a quick example:

Consider the sentence: "The quick brown fox jumps over the lazy dog."

A simple word-based approach would have these nine words as the units. However, a tokenizer might break it down like this (the exact tokenization depends on the specific tokenizer):

```
["The", "quick", "brown", "fox", "jumps", "over", "the", "lazy", "dog",
"."]
```

(Here, each word and the period are separate tokens).

Another tokenizer might break "jumps" into "jump" and "s", or "lazy" into "la" and "zy".

**Security Relevance:** How does tokenization relate to security?

- **Bypassing Filters:** Malicious actors might try to craft inputs that bypass security filters by manipulating tokenization. For example, they might insert subtle variations or misspellings that get tokenized in unexpected ways, potentially evading keyword detection.
- **Input Length Limits:** Since LLMs have context window limits measured in tokens, understanding tokenization is crucial for assessing how much malicious content can be fed into the model. A seemingly short malicious phrase might actually be a larger number of tokens.
- **Adversarial Attacks:** Researchers are exploring how adversarial perturbations at the token level can affect model behavior.

# Embeddings

We touched on this earlier, but it's worth reiterating. Embeddings are the numerical representations of tokens (and sometimes sequences of tokens). These vectors capture the semantic meaning of the tokens. Similar tokens will have embeddings that are close together in the high-dimensional embedding space. Understanding how these embeddings are learned and used is important because biases or vulnerabilities in the training data can be reflected in these embeddings.

As we briefly touched upon, embeddings are numerical representations of tokens (and sometimes sequences of tokens) in a high-dimensional space. Think of it as turning words or pieces of words into lists of numbers that capture their meaning.

The key idea behind embeddings is that words or tokens with similar meanings or that are used in similar contexts will have embeddings that are "close" to each other in this numerical space. The "closeness" is usually measured by the distance between these vectors (lists of numbers).

Imagine a map where countries that are geographically close are also close on the map. Word embeddings work similarly, but instead of geographical closeness, they represent semantic (meaning-based) closeness.

**How are these embedding numbers calculated?**

LLMs learn these embeddings during their pre-training phase by analyzing massive amounts of text. The model tries to predict the surrounding words for a given word. Through this process, it learns to associate words that appear together frequently and in similar contexts. The learned weights of the neural network layers responsible for generating these representations become the word embeddings.

Think of it like this: the model reads countless sentences and learns that "cat" is often found near "dog," "meow," "furry," etc. The embedding for "cat" will then be adjusted to be closer in the numerical space to the embeddings of these related words.

**Why are embeddings important for security?**

- **Semantic Similarity in Attacks:** Attackers might use semantically similar but syntactically different words to try and bypass keyword filters. Understanding how closely these words are embedded can help in designing more robust security measures. For example, "free gift" and "complimentary item" might have very close embeddings.
- **Bias Detection:** Embeddings can also capture biases present in the training data. For instance, if the training data disproportionately associates certain professions with certain genders, this bias might be reflected in the embeddings. Identifying and mitigating these biases is crucial for building fair and secure AI systems.
- **Understanding Model Interpretation:** Analyzing embeddings can sometimes provide insights into how the model "understands" language and the relationships between concepts. This can be helpful in identifying potential blind spots or vulnerabilities in the model's reasoning.

# Attention Mechanisms

As we discussed, these allow the model to weigh the importance of different tokens in the input sequence when processing each token. Different types of attention mechanisms exist, and understanding their nuances can be relevant when analyzing a model's behavior and potential weaknesses. For instance, if an attention mechanism is not implemented correctly, it might be exploited.

As we've discussed, attention mechanisms are a core component of the Transformer architecture and are what allow LLMs to understand the relationships between different parts of the input sequence.

Think back to our crowded party analogy. When you heard your name, your brain didn't process all the sounds equally. It focused on the specific frequencies and patterns that matched your name. Attention mechanisms in LLMs work similarly by allowing the model to focus on the most relevant parts of the input when processing each word (or token).

Here's a slightly more technical way to think about it:

For each token in the input sequence, the attention mechanism calculates a set of "attention weights" that indicate how much importance the model should give to every other token in the sequence (including itself) when processing that specific token.

These attention weights are determined by calculating a "similarity score" between the current token and all other tokens. Different types of attention mechanisms use different ways to calculate this similarity, but the underlying idea is to identify which tokens are most relevant to the current one in the context of the task.

Once these weights are calculated, they are used to create a weighted sum of the representations (embeddings) of all the tokens. This weighted sum provides a context-aware representation of the current token, taking into account the influence of the other relevant tokens in the sequence.

**Why are attention mechanisms important for security?**

- **Understanding Context in Prompts:** The attention mechanism is crucial for the LLM to understand the context of user prompts. If an attacker tries to inject malicious instructions within a seemingly harmless prompt, the attention mechanism will play a role in whether the model recognizes and acts upon those instructions.
- **Identifying Anomalous Inputs:** Deviations in attention patterns for malicious inputs compared to normal inputs might potentially be detectable and used as a security signal. Researchers are exploring ways to monitor these patterns.
- **Influence on Output Generation:** The attention weights directly influence which parts of the input the model focuses on when generating its output. Understanding these weights can be important in analyzing why a model might generate a harmful or biased response to a specific (potentially adversarial) input.
- **Vulnerabilities in Implementation:** Like any complex component, the implementation of attention mechanisms could potentially have vulnerabilities that attackers might try to exploit.

# Context Window

This refers to the amount of text that the LLM can consider when generating an output. Modern LLMs have a limited context window. If the relevant information is outside this window, the model might not be able to access it. The size of the context window can impact the model's ability to handle long inputs and maintain coherence, and it can also have security implications (e.g., how much malicious input can be effectively processed).

The context window of an LLM refers to the amount of text that the model can take into account when processing an input and generating an output. It's like the short-term memory of the model.

Think of it like reading a book. To understand what's happening on the current page, you might need to remember what happened on the previous few pages. The context window of an LLM defines how many "previous pages" (in the form of tokens) the model can effectively "read" and remember.

Modern LLMs have a limited context window. This limit is typically measured in the number of tokens. For example, a model might have a context window of 2,000 tokens, 32,000 tokens, or even more. Once the input sequence exceeds this limit, the model might start to "forget" earlier parts of the input.

**Why is the context window important for security?**

- **Handling Long Inputs:** Attackers might try to overwhelm the model with very long inputs containing malicious instructions embedded far back in the text, hoping that the model "forgets" the earlier parts or that security filters applied to shorter segments might miss them. Understanding the size of the context window is crucial for evaluating the effectiveness of such attacks.
- **Maintaining State in Conversations:** For conversational AI applications, the context window determines how much of the conversation history the model can remember. If the window is too small, the model might lose track of the conversation, which

could lead to security vulnerabilities if the conversation involves sensitive information or authentication steps.
- **Retrieval Augmented Generation (RAG):** Some systems enhance LLMs by allowing them to retrieve information from external knowledge bases to augment their context. The size and management of this extended "context" also have security implications. For example, if the retrieval mechanism is compromised, malicious information could be injected into the model's context.
- **Prompt Injection Attacks:** The effectiveness of prompt injection attacks can sometimes depend on where the malicious instructions are placed within the context window relative to the legitimate instructions.

The size of the context window is a significant factor in the capabilities and limitations of an LLM. Researchers are constantly working on increasing the context window size while maintaining efficiency and performance.

# Agentic AI

**Agentic AI** refers to the development of AI systems that can perceive their environment, make decisions, and take actions to achieve specific goals. Think of them as intelligent agents that can operate with some level of autonomy.

LLMs play a significant role in building these agents because of their ability to understand natural language, reason, and generate coherent text. This allows them to:

- **Perceive:** Understand instructions and information provided in natural language.
- **Plan:** Break down complex goals into smaller, manageable steps using their reasoning abilities.
- **Act:** Generate text-based commands or interact with other tools and systems through APIs (Application Programming Interfaces) to execute their plans.

Here are some key components often found in Agentic AI systems built with LLMs:

- **Planning Module:** The LLM itself often acts as the planner, taking a high-level goal and decomposing it into a sequence of actions.
- **Memory:** Agents often need to remember past interactions and information. This can be implemented through the LLM's context window (for short-term memory) or by integrating external memory systems (for longer-term memory).
- **Tool Use:** Agents can be equipped with the ability to use external tools, such as search engines, calculators, code interpreters, or even other specialized AI models, to gather information and perform specific tasks.
- **Execution Engine:** This component is responsible for carrying out the actions planned by the agent, often by interacting with external systems or generating outputs.
- **Observation Module:** After taking an action, the agent needs to observe the result and use that feedback to refine its plan and continue working towards its goal.

**Examples of Agentic AI Applications:**

- **Smart Assistants:** More advanced virtual assistants that can proactively perform tasks beyond just answering questions.
- **Autonomous Research Agents:** Systems that can conduct research on a topic, gather information from various sources, and synthesize findings.
- **Code Generation and Debugging Agents:** AI that can not only write code but also test and debug it autonomously.
- **Content Creation Agents:** Systems that can generate various forms of content, such as articles, social media posts, or even creative writing, with minimal human intervention.

**Security Considerations for Agentic AI:**

The autonomous nature of AI agents introduces new security challenges:

- **Unintended Actions:** If not carefully designed and controlled, agents might take actions that were not intended or have unforeseen consequences.
- **Vulnerability to Manipulation:** Malicious actors might try to manipulate the agent's planning or decision-making processes to achieve their own goals.
- **Data Security and Privacy:** Agents that interact with various systems and data sources need robust security measures to protect sensitive information.
- **Accountability:** Determining responsibility when an autonomous agent makes an error or causes harm can be complex.

Understanding how LLMs are used to build these agents and the potential security risks is crucial for your role in reviewing GenAI solutions.