# Intermediate ChatGPT

Attention scores: Numeric values computed during self-attention that indicate how much one token should attend to or weigh another token when forming contextual representations

ChatGPT: A conversational AI application built on large language models that generates human-like text responses to user prompts and can be customized with instructions or specialized GPTs

Custom GPTs: User-created, task-specific versions of GPT that combine tailored instructions, uploaded knowledge bases, and configuration settings to serve particular use cases or domains

Custom instructions: Persistent user-specified preferences provided to ChatGPT so it can personalize responses across all future conversations without repeating the same prompt details

Decoding: The process of converting a sequence of model-generated token IDs or probabilities back into human-readable text

Embedding: A numerical vector representation of a token that captures semantic and syntactic relationships and is used as input to neural networks

Few-shot prompting: A prompting technique that provides multiple examples within the prompt so the model can infer the desired format or reasoning pattern before producing an answer

Hallucination: When a model produces confident but incorrect or fabricated information, often due to gaps or biases in its training data or reasoning process

Hallucination: When a model produces confident but incorrect or fabricated information, often due to gaps or biases in its training data or reasoning process

Hidden state: The internal memory representation in recurrent models (like RNNs) that summarizes past inputs for use in future predictions

Large language model (LLM): A neural network trained on massive text datasets to predict and generate language patterns, typically by estimating the next token in a sequence

Masked-language modeling: A training task where certain tokens in a text are hidden (masked) and the model learns to predict those masked tokens from surrounding context

Multimodality: The capability of a model to process and generate multiple types of data (e.g., text, images, audio) so it can understand or respond across different input modalities

Next-token prediction: The training objective where a model learns to predict the most likely next token in a sequence given prior tokens, forming the basis of many generative language models

Parameters (weights): The learned numeric values inside a model (often millions or billions) that determine how inputs are transformed into outputs, analogous to synapses in a brain

PFSET writing structure: A five-part prompt framework (Persona, Framework, Specifications, Example, Topic) used to give an LLM clear role, structure, constraints, a sample, and subject context for better writing outcomes

Pre-training: The initial phase of training a model on broad, general-purpose text data so it learns language patterns and representations before task-specific tuning

Prompt engineering: The practice of designing and refining input prompts to guide a language model toward more accurate, relevant, or useful outputs

Recurrent neural network (RNN): A class of neural network that processes sequences step-by-step while maintaining a hidden state that carries information forward across time steps

Reinforcement learning from human feedback (RLHF): A training approach that uses human comparisons or preferences to rank model outputs and then optimizes the model to produce more preferred responses

Self-attention: A mechanism within transformers that lets each token weigh and aggregate information from every other token in the input sequence to build contextualized representations

Tokenization: The process of breaking text into discrete units (tokens) such as words, subwords, or characters so they can be mapped to numerical inputs for a model

Transformer architecture: A neural network design that processes all tokens in parallel and uses attention mechanisms to model relationships between tokens, enabling efficient handling of long sequences

Vanishing gradient problem: A training difficulty in deep or recurrent networks where gradients shrink during backpropagation, making it hard for the model to learn long-range dependencies

XML tags (in prompts): Structured markup inserted into prompts to delineate sections, roles, or reasoning steps so a model can more reliably parse and follow complex instructions