

1. GIVEN A SENTENCE WITH A MISSING WORD, APPLY A LANGUAGE MODEL TO PREDICT THE NEXT WORD IN SEQUENCE

What is Next Word Prediction?

Next word prediction means guessing what word should come next in a sentence. For example, if the sentence is "I am going to the," a model may guess the next word is "market." This is a main job of smart language models like ChatGPT.

How Does the Model Do It?

- **Breaking Text into Tokens:** The sentence is first broken into small parts called tokens. Tokens can be whole words or pieces of words or punctuation.
- **Turning Tokens into Numbers:** These tokens are turned into numbers so the model can understand them. Extra information is added to keep track of the order of words (this is called positional encoding).
- **Going Through the Model:** The tokens go through many layers inside the model. The model looks at the words before the missing word and uses a trick called self-attention to decide which earlier words are important.
- **Looking Only at Past Words:** When predicting a word, the model only looks at the words before it, not after. This is done using causal masking. It prevents the model from "cheating" by looking ahead.
- **Choosing the Next Word:** The model gives a score to every possible next word. These scores are turned into probabilities using something called the SoftMax function. The word with the highest chance is chosen.

What Happens After One Word?

If more words are needed, the new word is added to the sentence and the process starts again. This happens one word at a time until the sentence is complete.

Why is This Important?

Old methods like n-grams could only guess the next word using a few words before it. New models like transformers can remember and understand much more, making their predictions smarter and more natural.

Step-by-Step Example

- Input Sentence (with a missing word):** "I am going to the ____ to buy some groceries."
- Goal:** Use a language model to predict the most likely word that fits in the blank.
- Approach:** We use a pre-trained language model (like GPT) to fill in the blank based on context.

Prediction using a Language Model: Prediction: "market" So the completed sentence becomes: "I am going to the market to buy some groceries."

2. CLASSIFY THE TYPES OF TRANSFORMER MODELS. WITH A NEAT DIAGRAM ILLUSTRATE THE EVOLUTION OF THE TRANSFORMERS

Types of Transformer Models

- Seq2Seq (Encoder-Decoder)**
 - Uses both encoder and decoder.
 - Good for translation and summarization.
 - Example: DALL·E 1.
- Encoder-only**
 - Uses just the encoder.
 - Best for understanding tasks like classification.
 - Example: BERT (uses masked language modeling).
- Decoder-only**
 - Uses only the decoder.
 - Great for text generation (one word at a time).
 - Example: GPT series (GPT-1 to GPT-4).

Evolution Highlights

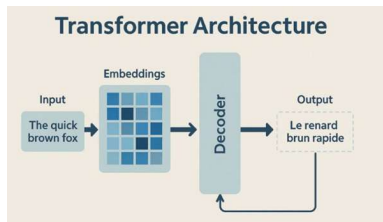


Figure 3.3: From the original transformer to GPT-4

- **2017 – Transformer:** Self-attention + encoder-decoder model.
- **2018 – BERT:** Encoder-only with bidirectional context.
- **2018 – GPT-1:** Decoder-only, left-to-right generation.
- **2019 – GPT-2:** Bigger, better generation.
- **2020 – GPT-3:** Few-shot, zero-shot learning at scale.
- **2021 – InstructGPT:** Tuned to follow instructions.
- **2023 – GPT-4:** Multimodal (text + images).
- **2023 – LLaMA 2:** Efficient and accessible.
- **2023 – Claude 2:** Better memory and safe responses.

3. WITH A NEAT DIAGRAM ILLUSTRATE THE ARCHITECTURE OF A TRANSFORMER

Transformer Architecture: A Game-Changer in NLP



Core Structure

a. Encoder Stack

- Converts input tokens into embeddings.
- Adds positional encoding to retain word order.
- Each layer has:
 - Multi-Head Self-Attention (MHSA): Lets tokens attend to each other for context.
 - Feed-Forward Network (FFN): Refines token representations.
- Outputs a context-rich representation ("memory").

b. Decoder Stack

- Generates output sequence one token at a time.
- Takes previous tokens and encoder output as input.
- Each layer has:
 - Masked MHSA: Prevents peeking at future tokens.
 - Encoder-Decoder Attention: Focuses on encoder output.
 - FFN: Further refines output.

c. Output Layer

- Decoder output → Linear layer → SoftMax → Next word prediction.

4. WITH A SUITABLE EXAMPLE ILLUSTRATE THE PROMPT ELEMENTS AND THE STRUCTURE

What Is a Prompt?

A prompt is an instruction you give to a language model. It includes what you want the model to do, the input data it should use, and how the response should be structured. It guides the model toward producing the desired output.

Key Elements of a Good Zero-Shot Prompt

- **Instruction**
 - **Purpose:** Tells the model what task to perform
 - **Example:** "Summarize the following text in one sentence."
- **Context**
 - **Purpose:** Provides background or useful information (optional)
 - **Example:** "The text is about the benefits of renewable energy."
- **Input**
 - **Purpose:** The actual content or data to be processed
 - **Example:** "Renewable energy like solar and wind reduces emissions..."
- **Output Cue**
 - **Purpose:** Suggests the expected format or start of the output
 - **Example:** Starts with: "Renewable energy sources, such as..."

Example Prompt Structure

Summarize the following text in one sentence.

The text is about the benefits of renewable energy.

Input: Renewable energy sources like solar and wind reduce greenhouse gas emissions and help fight climate change.

Output: Renewable energy sources, such as

Example Model Output

Renewable energy sources, such as solar and wind, help reduce emissions and combat climate change.

Why Use This Structure?

- Effective for zero-shot tasks (no prior task-specific training needed).
- Clearly separates what the model should do, what data to use, and what the output should look like.
- Helps produce more accurate and structured responses.

5. EXPLAIN SITUATIONAL PROMPTING OR ROLE PLAY IN LLMS WITH A SUITABLE EXAMPLE

What Is Situational Prompting (Role-Play)?

Situational prompting, also called role-play, is a method where a language model adopts a specific role or character in a simulated scenario. Unlike simple task prompts, this technique allows the model to interact dynamically as the situation evolves, responding in a context-aware manner.

Key Characteristics

- **Role-Based Identity:** The model acts as a character with specific traits (e.g., teacher, doctor).
- **Simulated Scenario:** The model engages in a fictional or realistic environment (e.g., consultation).
- **Dynamic Responses:** The model adapts replies based on user input and evolving context.
- **Goal-Oriented Dialogues:** Interaction is usually aimed at achieving a specific objective.

Example: Virtual Personal Stylist for StyleSprint

- **Scenario:** StyleSprint (a fictional fashion brand) wants to offer a personal stylist via chat.
- **Initial Prompt:** "Hi! I'm your personal stylist for today. What's the occasion you're dressing for?"
- **User Response:** "I'm attending a summer wedding."

- **Follow-up by LLM:** "Do you prefer bold colors or pastel shades?"
- **User Input:** "I like pastel shades."
- **Final Recommendation:** "For a summer wedding, I recommend our Floral Maxi Dress paired with the Vintage Sun Hat. It's elegant and perfect for an outdoor setting!"

Why Use Situational Prompting?

- Creates interactive experiences beyond static Q&A.
- Allows the model to personalize advice or services.
- Enhances engagement and realism in simulations.
- Useful for education, customer service, gaming, healthcare, and more.

This method blends persona with context-driven interaction, enabling language models to act like virtual assistants, coaches, consultants, or companions—adapting as the conversation unfolds.

6. DEFINE NEURAL NETWORK LANGUAGE MODEL. EXPLAIN NNLM WITH A SUITABLE EXAMPLE

What Is a Neural Network Language Model (NNLM)?

The NNLM, introduced by Yoshua Bengio et al. in 2003, was a breakthrough in natural language processing. It was the first neural approach designed to predict the next word in a sentence using learned word embeddings and hidden layers.

Key Features of NNLM

- **Goal:** Predict the next word in a sequence (language modeling).
- **Input:** A fixed-size context of previous words (e.g., "The cat is").
- **Architecture:** Includes an embedding layer, hidden layer, and output layer.
- **Word Embeddings:** Learns compact vector representations of words capturing semantic meaning.
- **Advantage Over Count-Based Models:** Generalizes better to unseen word combinations using continuous space.
- **Limitations:** Struggled with longer sequences and large vocabularies due to memory and compute constraints.

Simplified Example

Let's say the model is trained on the sentence:

"The cat is sleeping"

Suppose we want to predict the next word after "The cat is".

Step-by-Step (Conceptual)

- Input Words:** "The", "cat", "is"
- Convert to Embeddings:**
 - "The" → [0.2, 0.4, 0.1]
 - "cat" → [0.6, 0.1, 0.7]
 - "is" → [0.3, 0.5, 0.2]
- Feed to Hidden Layer:** The model combines these vectors and processes them in a hidden layer.
- Output Layer (Softmax):** Predicts probability distribution over vocabulary.
- Predicted Word:** "sleeping" (if it has the highest probability).

7. ANALYZE WHY TRANSFORMERS HAVE LARGELY REPLACED RNN-BASED MODELS IN ADVANCED NLP TASKS. WHAT ARE THE ARCHITECTURAL ADVANTAGES?

Why Transformers Replaced RNNs in NLP

- Limitations of RNNs and LSTMs**
 - RNNs process data step-by-step (sequentially), passing information along the sequence.
 - They have trouble learning from very long sequences due to the vanishing gradient problem (early inputs get "forgotten" during training).
 - LSTMs improve on RNNs by using gates to keep important information over longer distances.
 - However, LSTMs still process input sequentially, which slows down training and inference.
 - This sequential nature limits their efficiency on long texts and parallel computation hardware (like GPUs).
- Transformer Architecture Advantages**
 - Introduced in 2017, Transformers process entire sequences in parallel rather than step-by-step.
 - The core innovation is the self-attention mechanism, which lets the model weigh the importance of every word relative to all others in the sentence, no matter their distance apart.
 - Multi-head attention means the model looks at different aspects of the input simultaneously, capturing richer relationships.
 - Since attention does not encode order naturally, positional encoding is added to keep track of word positions.
 - The Transformer has two main parts:
 - **Encoder:** Processes input sequences into meaningful representations using self-attention and feed-forward layers.
 - **Decoder:** Generates output sequences, using masked self-attention (to avoid "seeing" future words) and attends to encoder outputs.
 - This parallel processing enables better handling of long-range dependencies (context from far apart words).
 - Transformers are scalable to very large models with billions of parameters, enabling state-of-the-art performance.

8. ANALYZE WHY LLMs MAY PRODUCE BIASED OR FAULTY INCORRECT OUTPUTS. EXPLAIN THE MAJOR GUIDELINES TO BE FOLLOWED WHILE USING LLMs

Why LLMs Produce Biased or Incorrect Outputs:

- **Hallucination:** LLMs sometimes generate plausible-sounding but factually false information because they predict text based on patterns, not real-world facts or logic.
- **Bias and Toxicity:** Training on large, unfiltered datasets causes LLMs to learn societal biases (e.g., gender, racial) and reproduce or amplify stereotypes and toxic language.
- **Static Knowledge:** LLMs only know what was in their training data up to a cutoff date; they cannot update themselves with new information or events in real time.
- **Accidental Memorization:** Models may memorize and regurgitate rare or sensitive data from training, risking privacy leaks or exposure of confidential info.

- **Ambiguous Prompts:** Vague or unclear inputs cause LLMs to “guess” user intent, potentially leading to inaccurate or unintended responses.
- **Susceptibility to Jailbreaking:** Attackers can exploit weaknesses or use tricks to bypass content filters, causing LLMs to generate restricted or harmful content.

Best Practices for Responsible LLM Use:

- **Understand Limitations:** Know that LLMs don’t reason or verify facts; they generate based on learned patterns, which means errors and bias are inherent.
- **Apply Domain Expertise:** Use subject-matter knowledge to critically evaluate outputs, especially in fields requiring precision like medicine, law, or coding.
- **Verify and Avoid High-Stakes Use Without Validation:** Always check model outputs before acting on them. Avoid relying on LLMs alone in critical scenarios where mistakes could be harmful.
- **Address Bias and Toxicity:** Use curated training data, fine-tune models carefully, apply post-processing filters, and include human review to mitigate biased or toxic outputs.
- **Give Clear, Specific Prompts:** Provide unambiguous, detailed instructions to reduce model guesswork and improve output accuracy.
- **Use Grounding Methods:** Enhance factuality by combining LLMs with verified external data sources (e.g., Retrieval Augmented Generation) to anchor responses in truth.
- **Consider Ethical Responsibilities:** Promote transparency, fairness, privacy, accountability, and inclusivity in model design and deployment.
- **Implement Robust Safeguards:** Use filtering, safety fine-tuning, continuous monitoring, and safeguards against jailbreaking to prevent misuse or harmful generation.

9. EXPLAIN HOW RECURRENT NEURAL NETWORKS (RNNs) FUNCTION IN THE CONTEXT OF LANGUAGE MODELLING. WHAT ARE THEIR STRENGTHS AND WEAKNESSES COMPARED TO TRADITIONAL LANGUAGE MODELS?

How RNNs Work in Language Modeling:

- RNNs are neural networks designed to handle sequences of data (like sentences).
- They process input one element at a time, keeping an internal state that stores information from previous steps.
- Unlike regular feedforward networks that treat each input independently, RNNs pass information along the sequence, capturing context.
- This allows RNNs to understand relationships over time or across words in text.
- LSTM networks are a special type of RNN made to solve the vanishing gradient problem, which limits learning from long sequences.
- LSTMs use gates to control information flow, enabling them to remember important information from both recent and distant parts of the sequence.

Strengths of RNNs (especially LSTMs):

- They work well when previous context affects future predictions.
- Can handle long sequences and complex dependencies in language.
- Achieved top performance in language tasks and text classification.
- Became the main neural architecture in NLP due to capturing both short- and long-range context.

Weaknesses of RNNs:

- Standard RNNs struggle with very long sequences due to the vanishing gradient problem.
- They process sequences step-by-step (sequentially), which makes training slower compared to models that can run in parallel.

Comparison to Traditional Language Models:

- Traditional NLP methods include n-grams, hidden Markov models (HMMs), count vectors, and TF-IDF.
- N-grams and HMMs capture local and short-range dependencies by looking at a fixed number of previous words.
- Count-based methods extract basic word frequencies but fail to understand word meanings or semantics in context.
- RNNs and LSTMs improved on this by:
 - Managing long-range dependencies beyond fixed window sizes.
 - Using word embeddings that capture semantic meanings, unlike count-based vectors.
 - Providing a richer, context-aware representation of language.

10. IDENTIFY AND DISCUSS DIFFERENT TYPES OF PROMPTS AND THEIR PURPOSES IN INTERACTION WITH LANGUAGE MODELS. HOW DO VARIATIONS IN PROMPT DESIGN AFFECT MODEL OUTPUTS?

What is Prompt Engineering?

- It’s the skill of designing instructions (prompts) that guide a language model’s responses without changing the model itself.
- Combines technical know-how, creativity, and understanding human psychology to get accurate, relevant, and desired outputs.
- Prompts act like directions to focus the model on the right task.

Types and Variations of Prompts:

a. Basic Prompt Elements

- **Instruction:** Clear command or question.
- **Context:** Background info needed to understand the task.
- **Input:** Data or content the model works on.
- **Output cue:** Signals the format or style of the answer (e.g., JSON, Markdown).
- Clear structure helps models respond more reliably and in the expected format.

b. Zero-Shot Prompting

- The model answers based only on the prompt and its pre-trained knowledge, with no examples.
- Responses can be generic or less precise without example guidance.

c. Few-Shot Prompting (In-Context Learning)

- Provides a few task examples inside the prompt.
- Enables the model to mimic examples, improving accuracy and consistency.
- Reduces the need for expensive fine-tuning.
- Success depends on example quality and relevance.
- Can convey specific styles or tones (like brand voice).

d. Behavioral Prompting (Elevating Prompts)

- Uses emotional cues, personas, or role-playing to influence tone and style.
- **Emotional cues:** Positive words can improve performance and engagement.
- **Personas:** Asking the model to adopt a character shifts responses toward specific viewpoints or tones.

- **Situational prompts/Role-play:** Creates dynamic, context-aware dialogues.
- e. Prompt Chaining**
 - Breaks complex tasks into multiple prompts, feeding outputs back as inputs.
 - Allows iterative refinement and better steering of responses over a conversation.
- f. Retrieval Augmented Generation (RAG)**
 - Augments prompts with relevant external information from databases or documents.
 - Grounds outputs in factual data, improving reliability and factual accuracy.
 - Overcomes context window limits by storing and retrieving external “memory.”
 - Can reduce harmful or biased responses by controlling the context given to the model.
 - Often combined with few-shot examples and chaining.

11. EXAMINE THE KEY GUIDING PRINCIPLES FOR INTERACTING EFFECTIVELY WITH LANGUAGE MODELS. HOW DO THESE PRINCIPLES ENHANCE THE OVERALL PERFORMANCE OF GENERATED RESPONSES?

- a. Understand LLMs’ Limitations**
 - LLMs can make mistakes and sometimes produce false but plausible-sounding info (called hallucination).
 - Knowing this helps users verify outputs and avoid blindly trusting them.
 - This leads to safer, more reliable use.
- b. Use Domain Knowledge**
 - Experts in the relevant field can spot errors and correct them.
 - This is crucial in areas needing precision, like coding, medicine, or research.
 - Expertise improves accuracy and trustworthiness of results.
- c. Ground Outputs with Facts**
 - Add reliable, factual info to the prompt or use retrieval methods like RAG to supply verified data.
 - Grounding reduces hallucination and guides the model to produce fact-based answers.
 - Combining retrieval with filtering helps prevent harmful or misleading responses.
- d. Avoid Ambiguity in Prompts**
 - Vague prompts cause the model to guess or infer intent, which may lead to wrong answers.
 - Clear, structured prompts with explicit instructions, context, inputs, and output cues guide the model well.
 - This clarity results in responses that better match user expectations.
- e. Don’t Use LLMs Without Verification or in High-Risk Cases**
 - Avoid relying on LLM output when you cannot check correctness.
 - Don’t use LLMs for critical decisions where mistakes have serious consequences.
 - This principle ensures responsible use and prevents harm.