

How LLMs Work & Why Prompt Engineering Matters



ByteMonk
204K subscribers

Join

Subscribe

88



Share

Thanks

Clip



839 views Jun 21, 2025 [Members first on June 19, 2025](#) [#GPT4](#) [#PromptEngineering](#) [#agentical](#)

Prompt engineering isn't about tricks — it's about understanding how large language models (LLMs) actually work.

In this first part of the series, we unpack the evolution of LLMs — from Seq2Seq to Transformers — and explain how models generate responses, why prompting replaced fine-tuning, and what it means to design prompts that shape output reliably. We also touch on the system design behind prompt-driven workflows, memory injection, and multi-modal prompts — setting the stage for building real-world AI applications. This is where solid architecture meets practical AI.

My LinkedIn Profile: [in / bytemonk](#)

★ Timestamps

0:00 – Intro: Why Prompt Engineering Matters

1:05 – What Is a Language Model Really Doing?

1:46 – Early Models: Seq2Seq and Thought Vectors

2:28 – Bottlenecks in Seq2Seq (and Why It Failed)

3:56 – Attention Mechanism and "Attention Is All You Need"

4:50 – Birth of Transformers (Parallelism, Power, Limitations)

5:39 – Enter GPT: From Fine-Tuning to Prompting

6:49 – Prompts vs Completions: The Heart of LLMs

8:20 – Predicting Continuations with Real-World Patterns

8:32 – Prompt Engineering as a System, Not Just a Prompt

9:04 – Levels of Prompting: Context, Memory, Tools

LLM

Predicting what comes next in a sentence

Better test cases

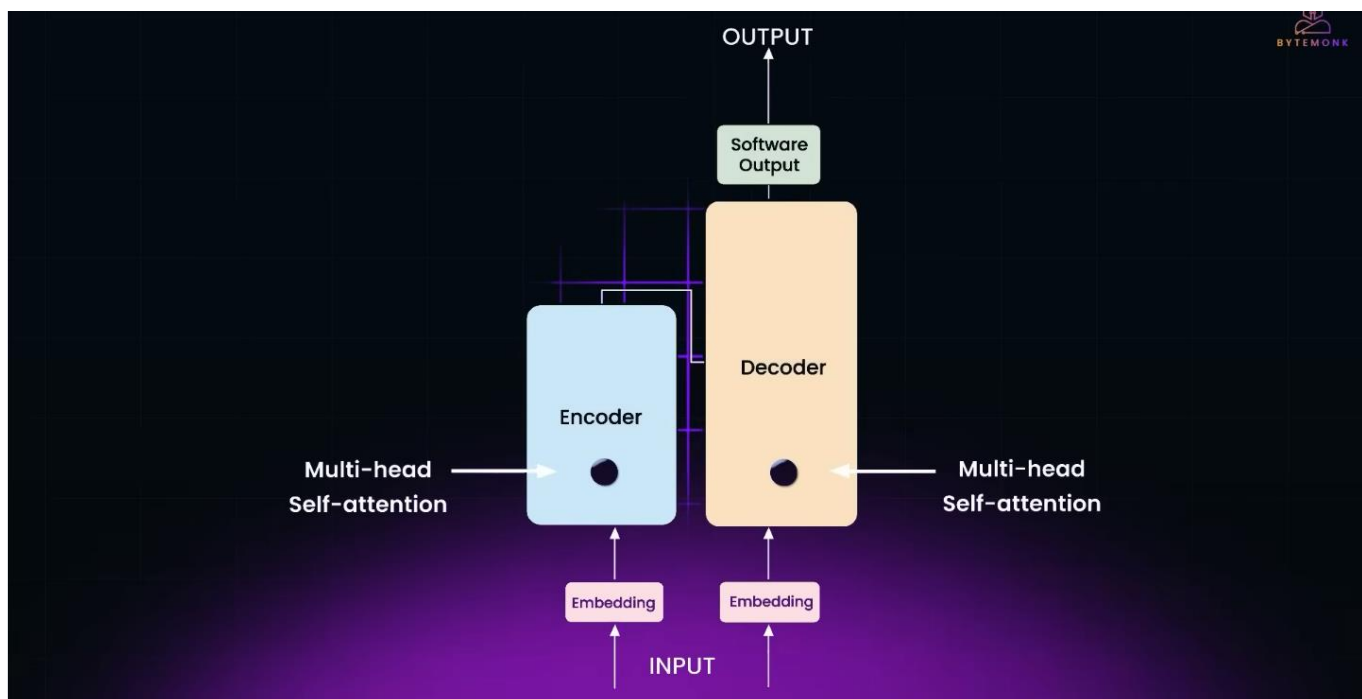
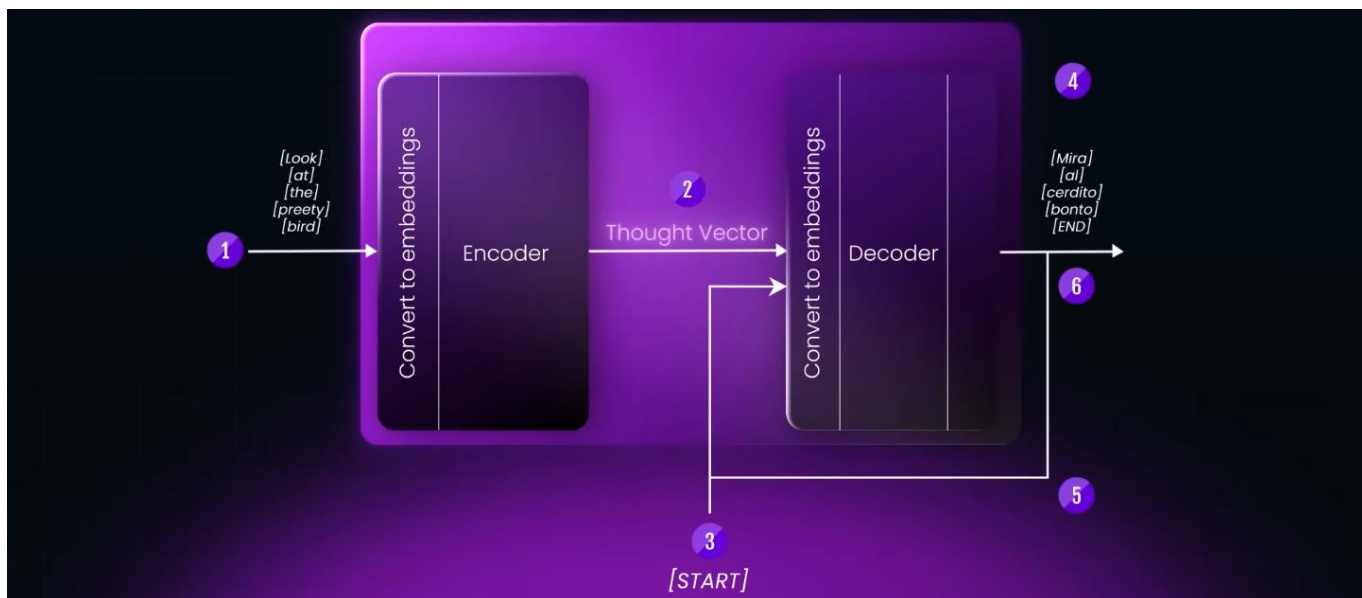
Boilerplate code

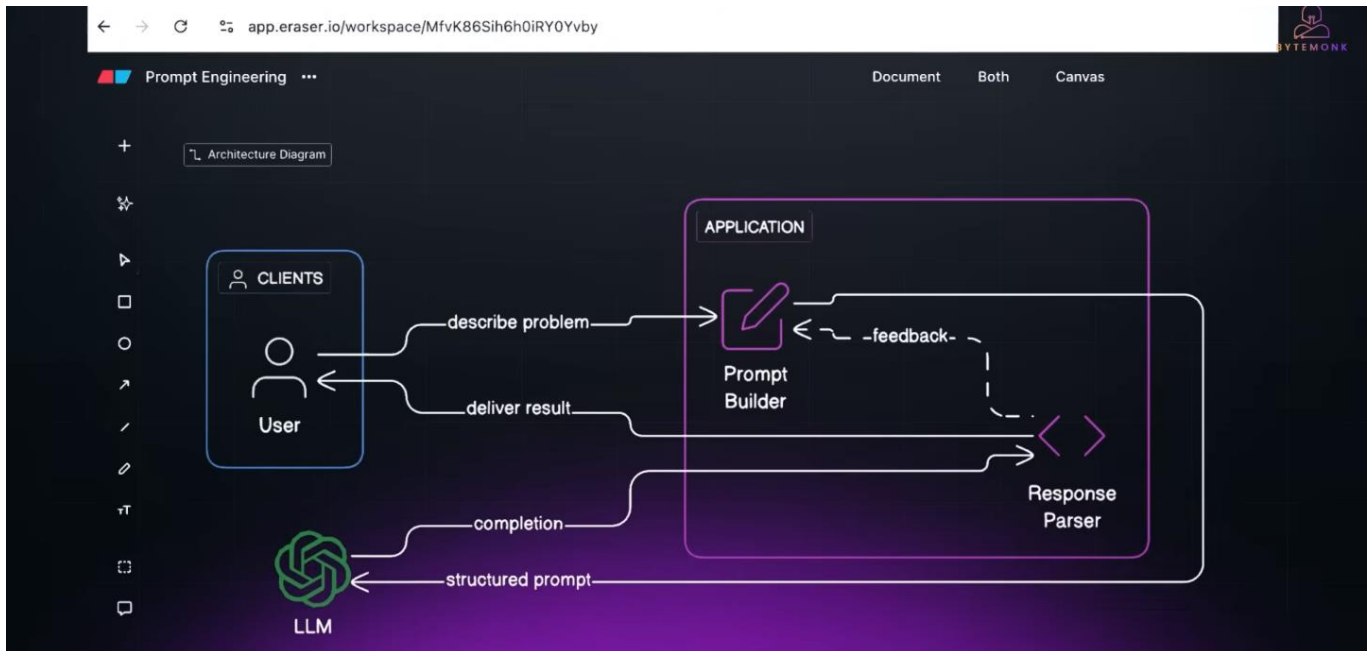
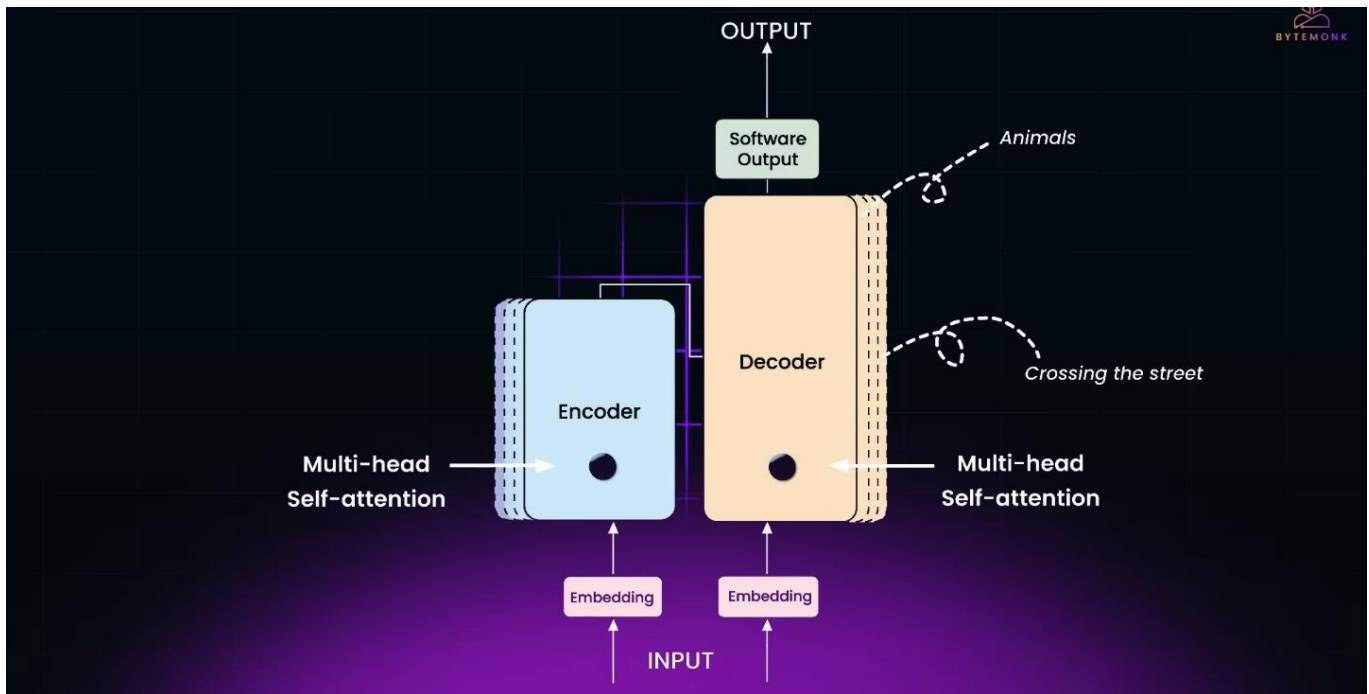
Debugging faster

Summarizing logs



Predicting the next word





LLM

Predict the next word in a sequence



Recurrent Neural Networks



[START]

These models processed 1 token at a time and updated an internal memory after each step.

Sentences - Paragraphs - Documents

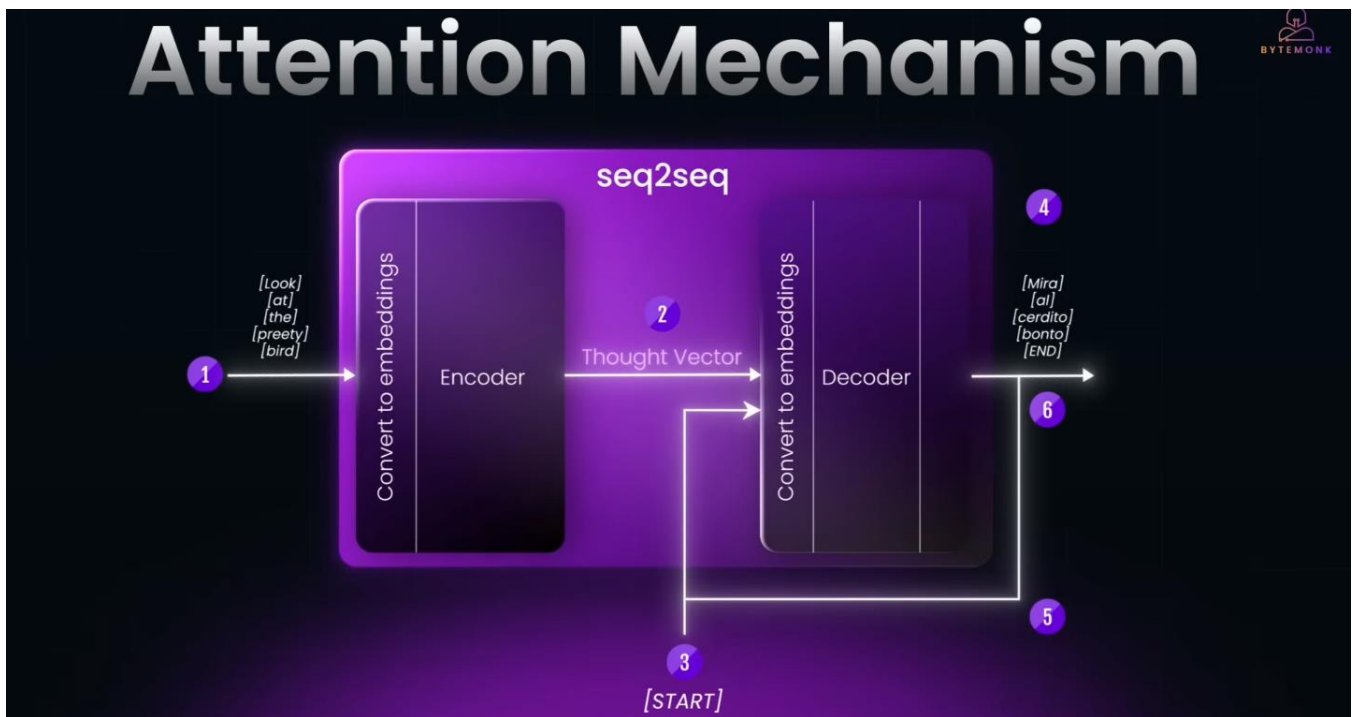
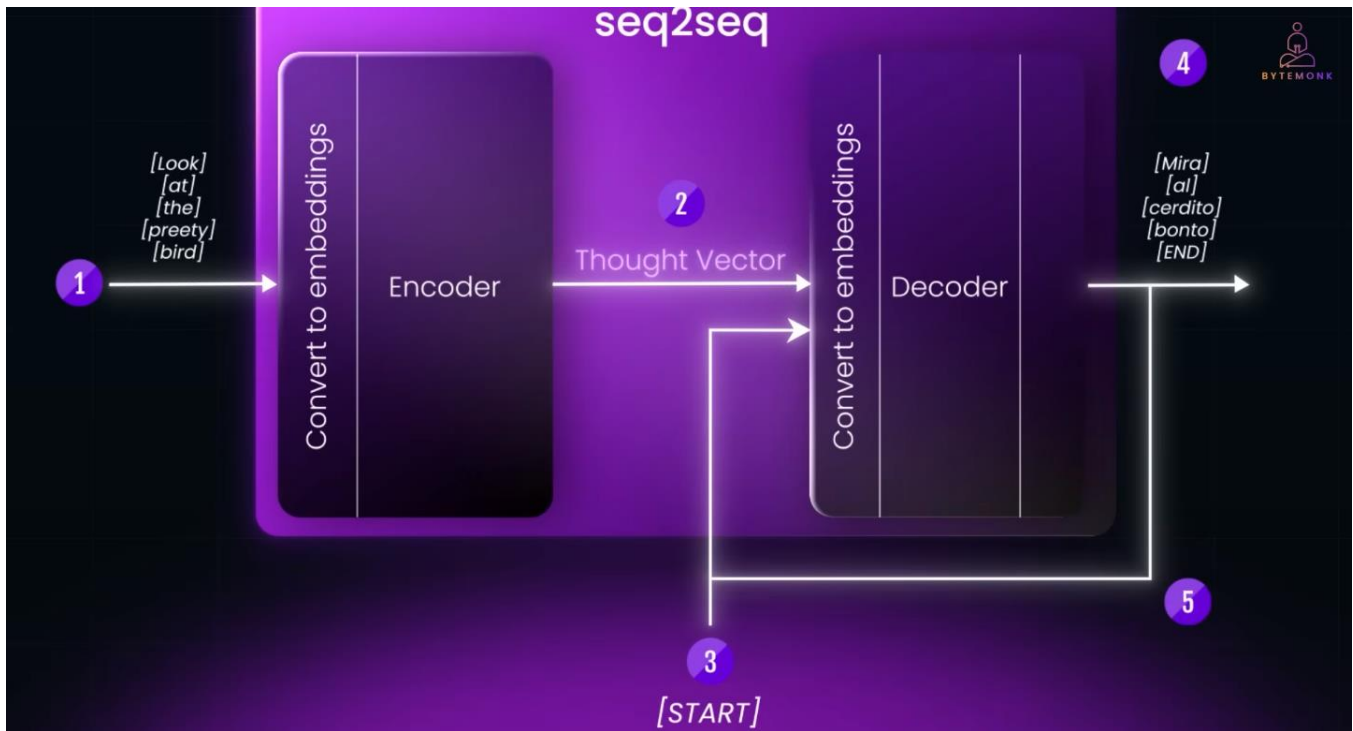


[START]

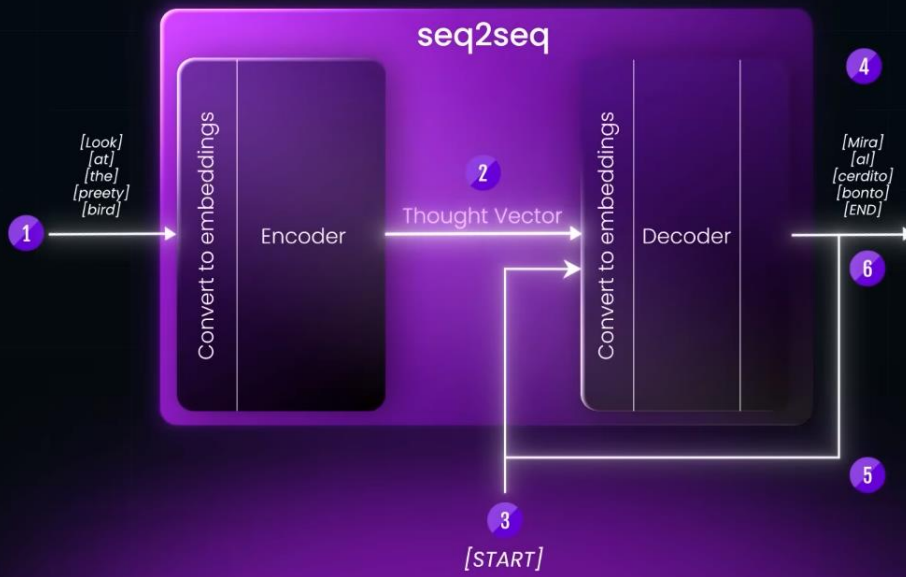
Information Bottleneck



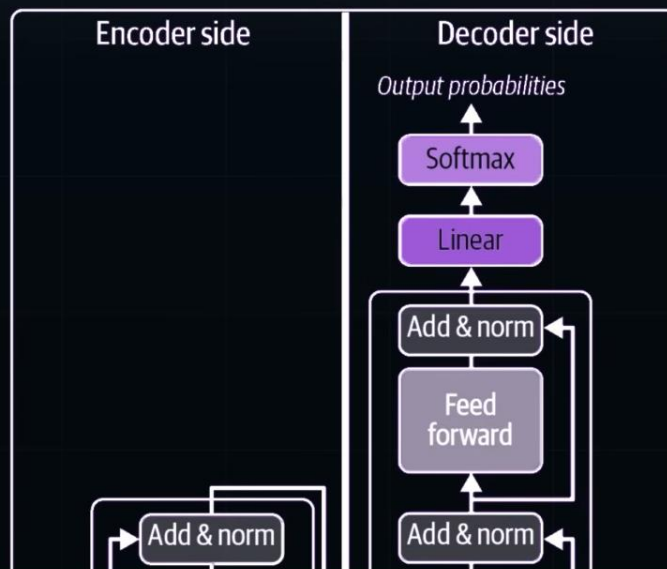
[START]



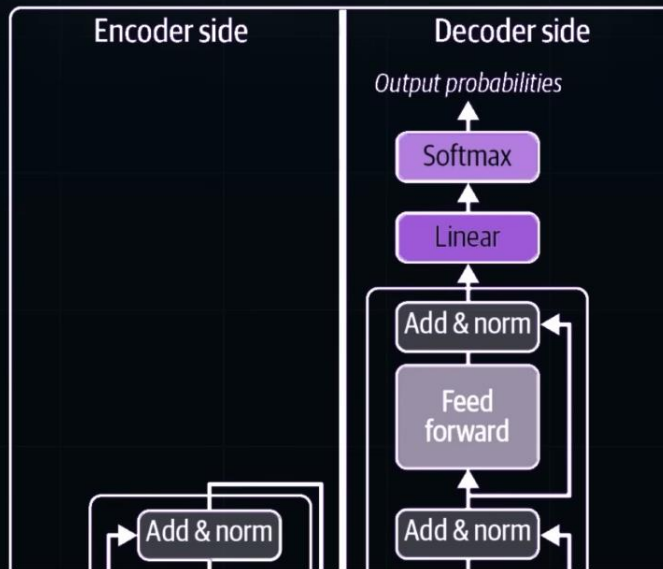
Last of those hidden states



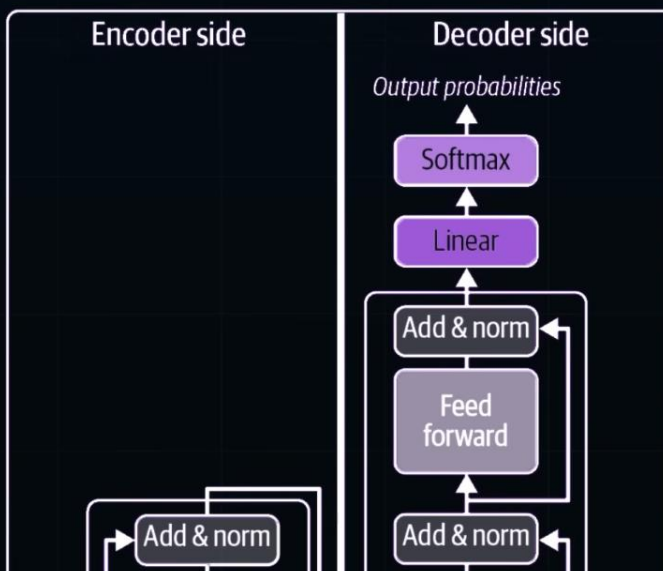
Keep all the hidden state vectors



Look back at all of them



FULL MAP of the input



Attention Is All You Need

Ashish Vaswani^{*}
Google Brain
avaswani@google.com

Noam Shazeer^{*}
Google Brain
noam@google.com

Niki Parmar^{*}
Google Research
nikip@google.com

Jakob Uszkoreit^{*}
Google Research
uszko@google.com

Llion Jones^{*}
Google Research
llion@google.com

Aidan N. Gomez^{*†}
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser^{*}
Google Brain
lukaszkaizer@google.com

Illia Polosukhin^{*}
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder-decoder configuration. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

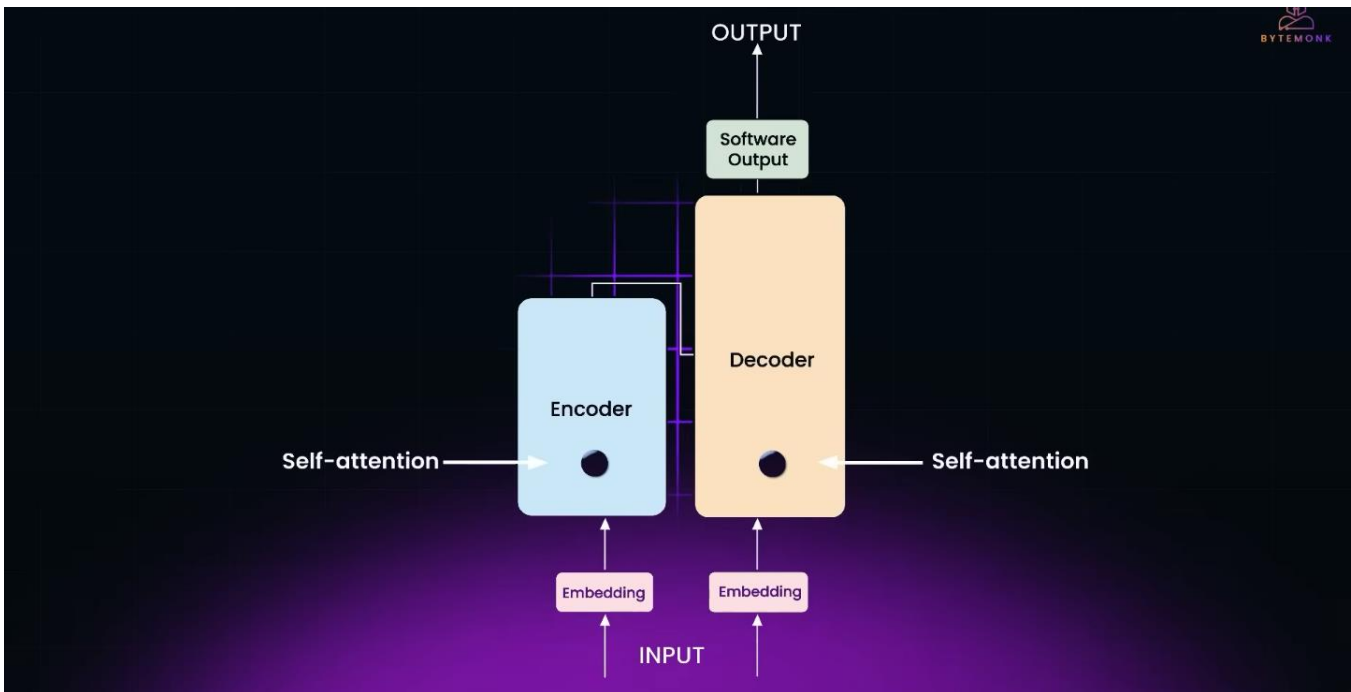
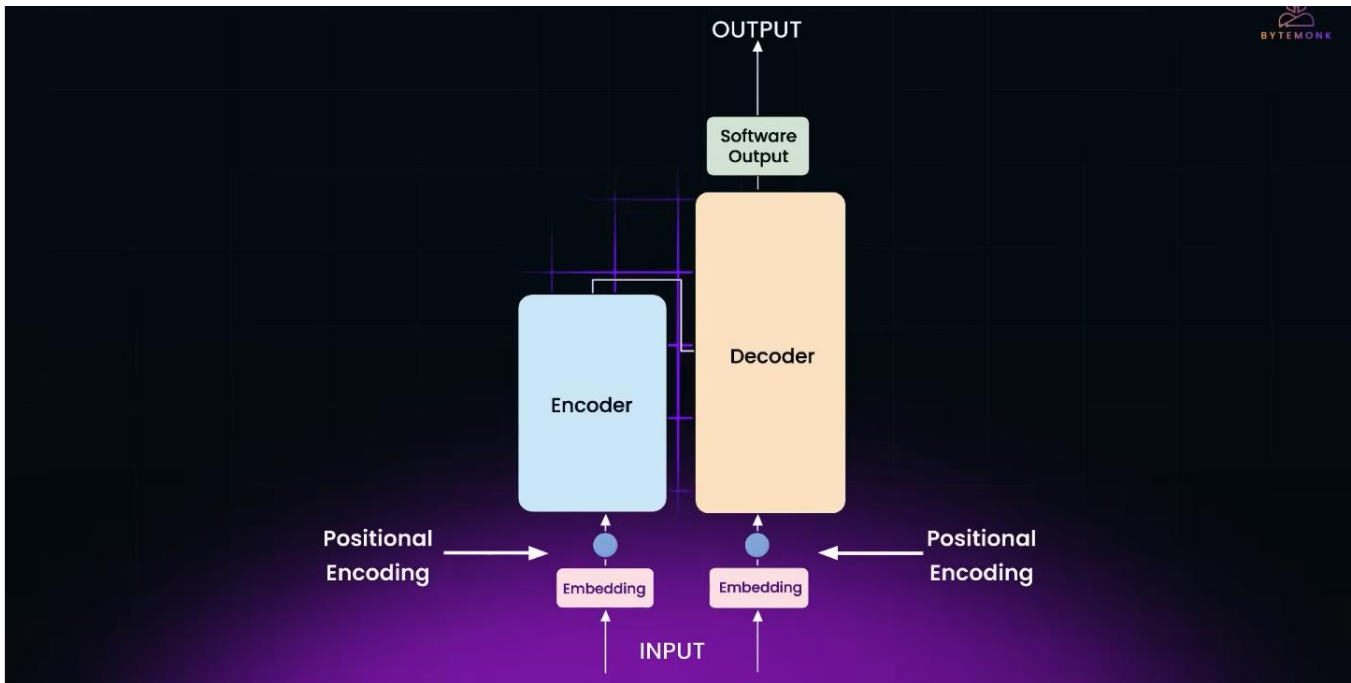
1 Introduction

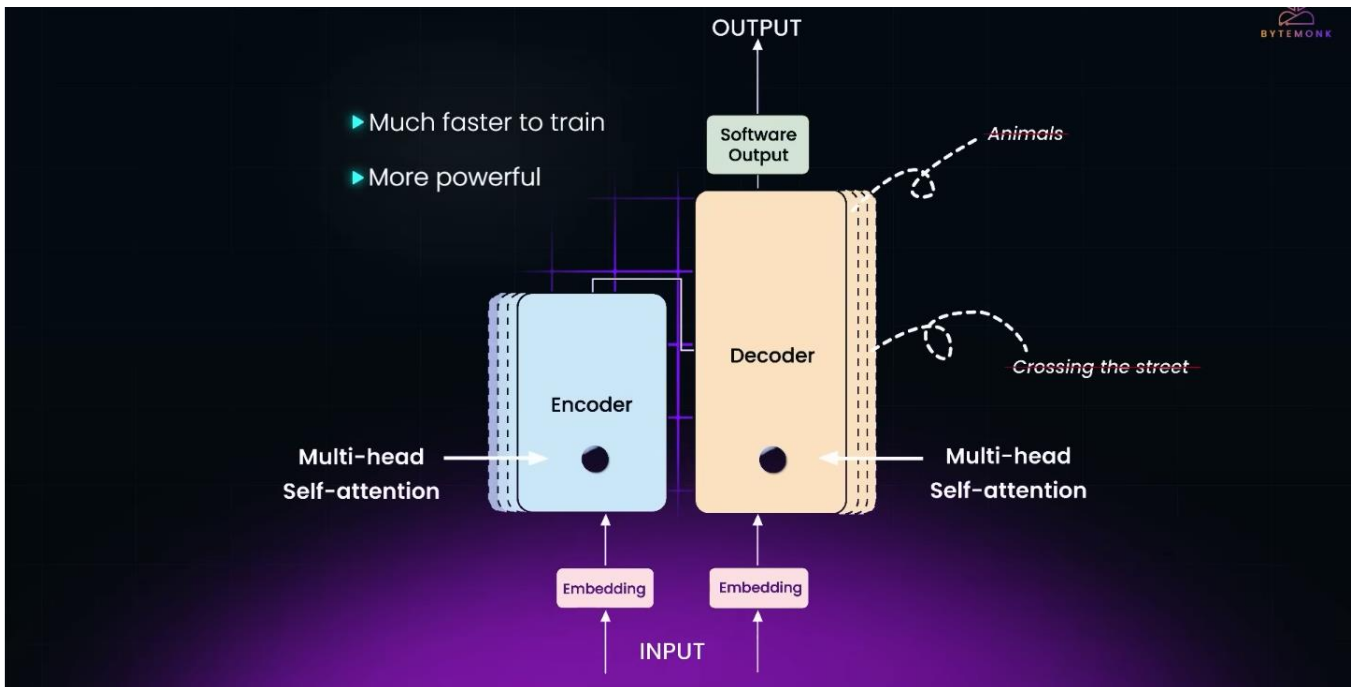
Recurrent neural networks, long short-term memory [12] and gated recurrent [7] neural networks in particular, have been firmly established as state of the art approaches in sequence modeling and transduction problems such as language modeling and machine translation [31, 2, 5]. Numerous efforts have since continued to push the boundaries of recurrent language models and encoder-decoder architectures [34, 22, 14].

Recurrent models typically factor computation along the symbol positions of the input and output sequences. Aligning the positions to steps in computation time, they generate a sequence of hidden states h_t , as a function of the previous hidden state h_{t-1} and the input for position t . This inherently sequential nature precludes parallelization within training examples, which becomes critical at longer

The diagram illustrates the internal structure of the Transformer's Encoder and Decoder. The **Encoder side** (left) consists of a stack of N identical layers. Each layer contains a **Multihead attention** sub-layer followed by a **Feed forward** sub-layer, with **Add & norm** operations and residual connections. The input is processed through **Input embedding** and **Positional encoding**. The **Decoder side** (right) also consists of a stack of N identical layers. Each layer contains a **Masked multihead attention** sub-layer, followed by a **Multihead attention** sub-layer (which takes encoder output as input), a **Feed forward** sub-layer, and **Add & norm** operations with residual connections. The output is processed through **Output embedding** and **Positional encoding** (shifted right). The final output is passed through a **Linear** layer and a **Softmax** layer to produce **Output probabilities**.

This high-level diagram shows the flow of data in an encoder-decoder model. **INPUT** is fed into an **Embedding** block, which then connects to the **Encoder** (blue box). The **Encoder** outputs to another **Embedding** block, which connects to the **Decoder** (orange box). The **Decoder** produces a **Software Output**, which is then fed back into an **Embedding** block and finally to the **OUTPUT**. A dashed line indicates a feedback loop from the **Software Output** back to the **Decoder's** input embedding.

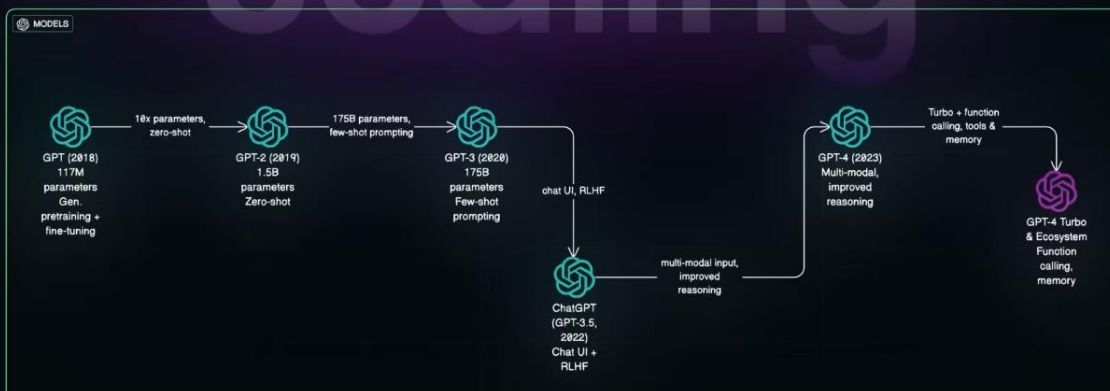


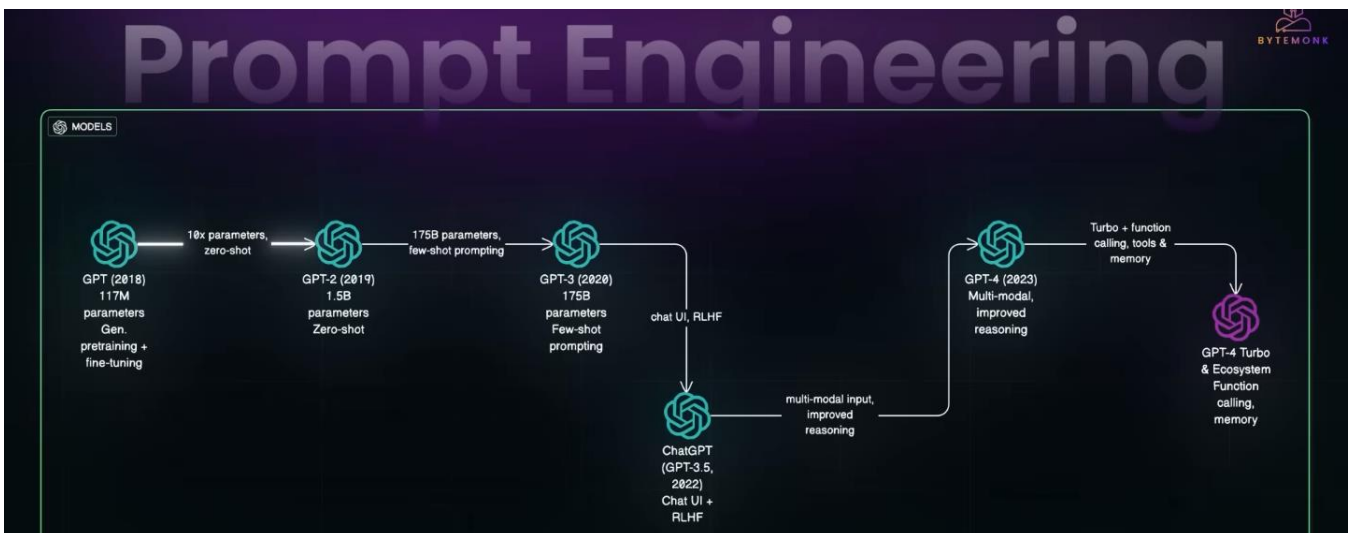
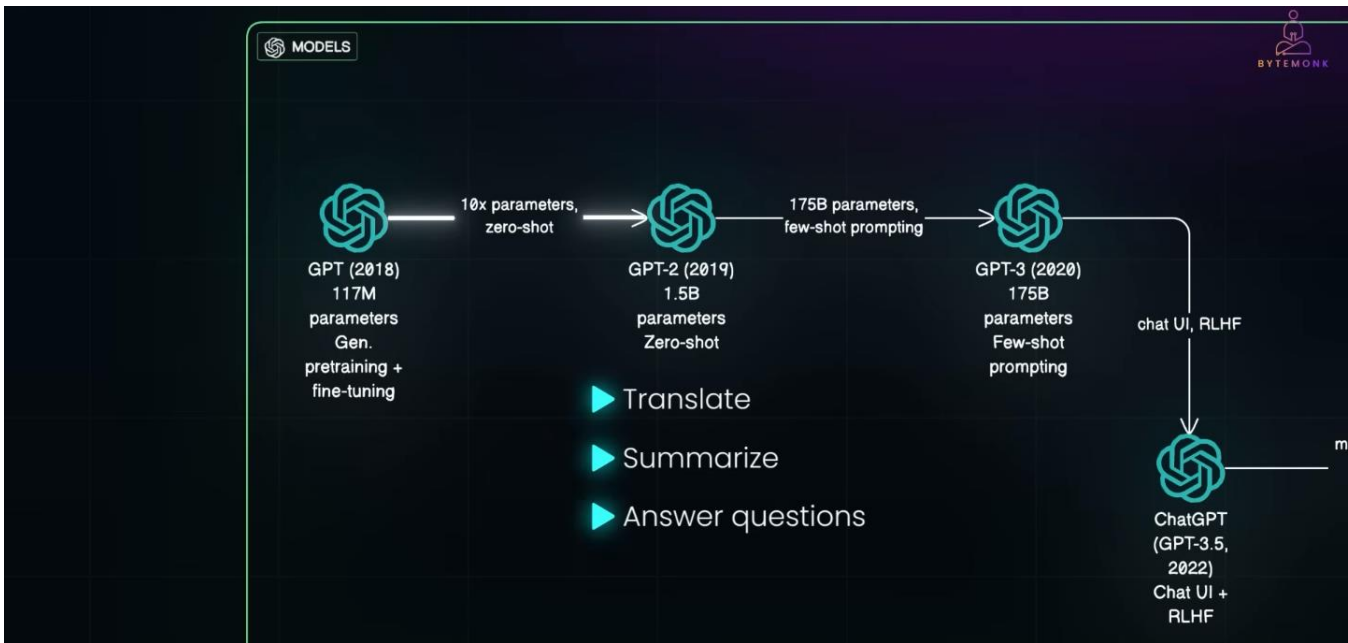


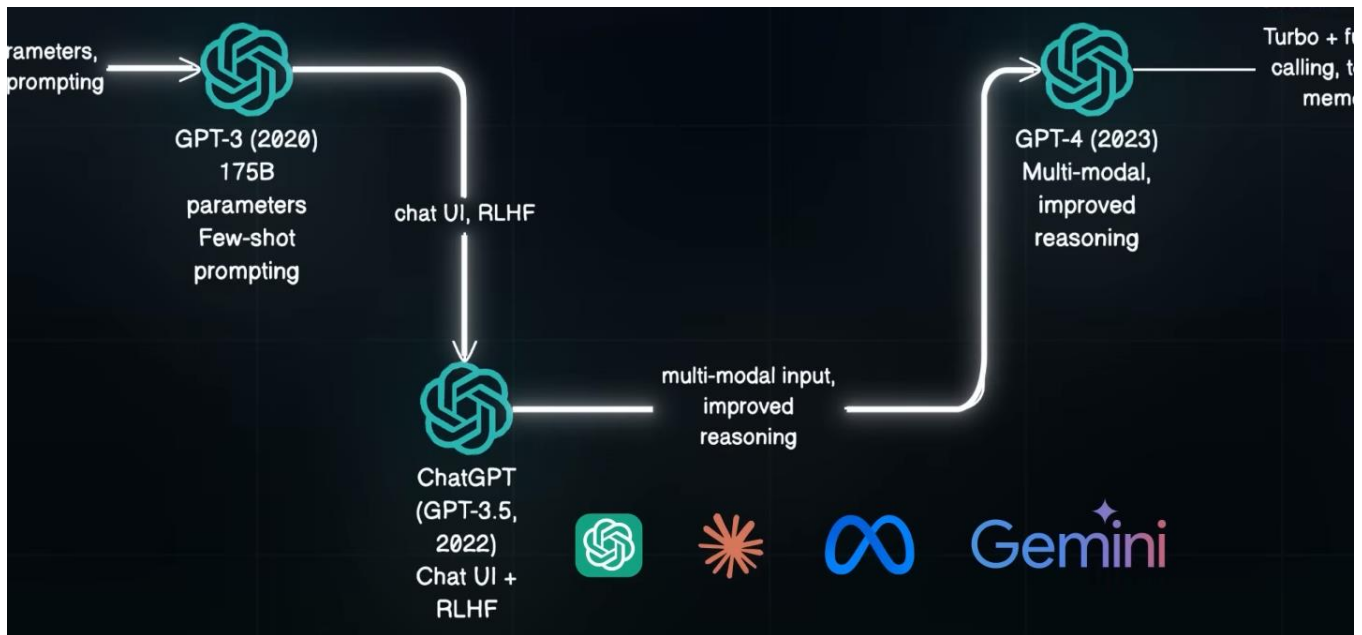
Tradeoff

Transformers can't handle unlimited input
Fixed context window

Scaling







String of text	Returns a continuation
Prompt	Completion

Prompt

The API request failed because...

Completion

The authentication token was missing

The deployment failed during staging.
Now the app won't start on...

g9dl93

Monday

Production

The deployment failed during staging.
Now the app won't start on production...

Let's refactor the pipeline before EOD.

Can you check the logs on your end?

Try rolling back to the last known stable build.

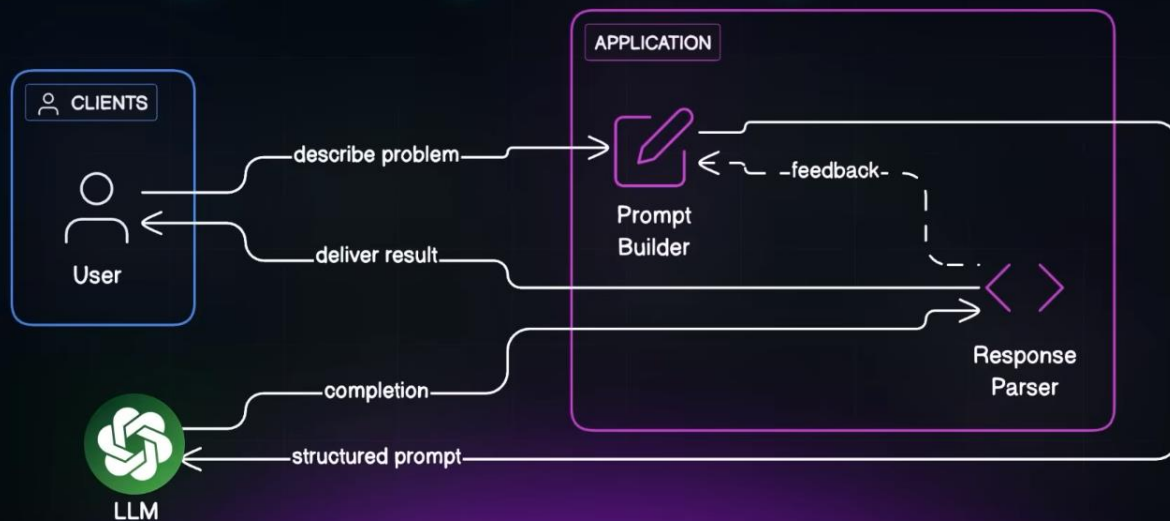
Data Patterns

Not logic or facts

▶ Smart

▶ Consistent

▶ Reliable



Basic level



Enriched Level

▶ Past Conversations

▶ Open browser tabs

▶ Search results

Enriched Level

Contextual
Memory Aware

Advance Level

▶ Call APIs

▶ Check calendars

▶ Send emails

Advance Level

- ▶ Manage tool use
- ▶ Handle errors
- ▶ Pass state
- ▶ Keep the loop safe

Prompt Engineering

