

Prompt Engineering for LLMs: A 2025 Guide

Fundamentals

Definition: Prompt engineering is the art of crafting and structuring input prompts to guide a large language model (LLM) toward desired outputs. As OpenAI explains, it is “the process of designing and optimizing input prompts to effectively guide a language model’s responses” ¹. By carefully choosing words, context, and format, we can align an LLM to our task without changing its weights ² ¹.

Importance: Effective prompt engineering is crucial for unlocking an LLM’s potential. The quality of your prompt directly influences the relevance, accuracy, and coherency of the model’s output ³. In practice, clear prompts can steer models to solve complex tasks (e.g. writing code or summarizing research) while vague prompts yield generic or irrelevant results ⁴ ³. In short, prompt engineering is about *steerability* and *alignment*: it lets humans shape LLM behavior through input alone, making these powerful AI tools practical for real-world use ² ³.

Core Principles

- **Clarity & Specificity:** Be unambiguous. Use precise language and explicit details so the model “knows exactly what to do” ⁵. For example, instead of “Tell me about robots,” ask “Summarize how autonomous robots can be used in warehouse automation.” Straightforward, specific wording anchors the model’s response to your intent ⁶ ⁵.
- **Context & Role Prompting:** Provide relevant context or assign a persona to shape style. Including background information or **system instructions** can anchor the response. For instance, a prompt might begin “You are a legal expert advising a startup...” to get formal, domain-specific answers. Role prompting (assigning a persona like “math tutor” or “marketing analyst”) is widely used to control style and depth ⁷. As LearnPrompting.org notes, instructing the model to adopt a role “shapes how it processes and delivers information” ⁷.
- **Step-by-Step Instructions (Chain-of-Thought):** Break tasks into steps or explicitly ask for reasoning. Chain-of-Thought (CoT) prompting, where you request the model to “think step by step,” often yields better reasoning on multi-step problems ⁸. For example, solving a math problem might be phrased “Explain your reasoning, then give the answer.” This guides the model to outline intermediate steps, improving accuracy on complex queries ⁸.
- **Examples & Demonstrations:** Use few-shot prompting by giving examples of inputs and outputs. Including exemplars of the desired format teaches the model what kind of response you want ⁹ ¹⁰. For example, to format answers as bullet lists, you might say: “List pros of X. Example: (‘Feature: Benefit’).” Examples reduce ambiguity and establish patterns that the model will follow ¹⁰ ⁹.

- **Iterative Refinement:** Expect to refine prompts. Prompt engineering is an iterative process: start with a draft prompt, evaluate the output, then tweak the wording, add context, or simplify instructions as needed ¹¹ ¹². As Palantir advises, use model feedback in a loop: “If the model misunderstands, adjust the wording and re-test” ¹². In practice, professionals treat the interaction like a conversation – multiple rounds of clarification often yield the best results.
- **Constraints & Formatting:** Explicit constraints help focus the response. Specify length, style, or forbidden content (e.g., “Summarize in three sentences,” or “List answers only, no extra text”). Palantir notes setting boundaries (“no more than three sentences”) and using negative instructions to avoid unwanted output ¹³. Clear constraints prevent the model from rambling or going off-topic.
- **Manage Length and Relevance:** Keep prompts concise and relevant to avoid overloading the model. Large prompts can hit context limits or confuse the model. Palantir’s guide recommends brevity (“provide necessary details without overloading”) and breaking complex tasks into parts ¹⁴. If your instruction is long or includes a lot of data (e.g. an article to summarize), consider summarizing or chunking the input first, or use retrieval augmentation to supply only the relevant parts.

Prompt Types

Prompting techniques can be categorized by how they guide the model:

- **Zero-Shot Prompts:** The simplest form: you give the task without examples. e.g. “*Translate this sentence to French: [text]*”. The model relies on its pre-training to perform the task. Zero-shot is quick but may be less accurate on novel tasks because the model hasn’t seen your specific format or intent ¹⁵.
- **Few-Shot Prompts:** You include a few input-output examples of the task before asking the model to continue. For example:

```
Q: "What is the sentiment of the sentence: 'I love this product!'"
A: "Positive."
Q: "What is the sentiment of the sentence: 'This was a bad experience.'"
A: "Negative."
Q: "What is the sentiment of the sentence: 'The movie was ok, not great.'"
A:
```

Few-shot prompts teach the model the desired format and criteria. As Lilian Weng notes, “few-shot learning often leads to better performance than zero-shot” because the model “first sees good examples” of the task ⁹.

- **Chain-of-Thought (CoT) Prompts:** You explicitly ask the model to show its reasoning or break down the steps. E.g. “Explain your reasoning step by step, then answer the question.” This encourages the model to generate intermediate “thoughts,” which improves multi-step problem solving ⁸. CoT is

especially effective on logic, math, and commonsense tasks. (Advanced CoT variants include Self-Ask, Tree-of-Thought, etc., which interleave retrieval or branch into multiple reasoning paths ⁸.)

- **Retrieval-Augmented Generation (RAG):** Here you combine the LLM with external knowledge sources. Instead of relying solely on the model's training data, you retrieve relevant documents or facts and feed them into the prompt (or use a tool call). For example, before asking a question, the system might fetch related text from a database. IBM explains that RAG “connects an LLM to a database and automates information retrieval to augment prompts with relevant data for greater accuracy” ¹⁶. Use RAG when your prompt needs up-to-date or domain-specific information beyond the model's internal knowledge. In practice, this may involve building a separate retrieval pipeline or using frameworks (LangChain, LlamaIndex) to fetch and insert context.
- **Function-Calling / Tool Use:** Modern LLM APIs (e.g. OpenAI's) allow the model to call external functions. In this approach, you provide a schema or function definitions to the model and instruct it to output a structured call when appropriate. For instance, you can give the model an “weather API” function and prompt: “Give me the weather for Paris.” The model then outputs a JSON with a function name and arguments, instead of a normal text reply. This makes the output machine-readable and lets developers hook the LLM into external systems. Function-calling is a powerful prompt type when you need structured data or to integrate LLMs with applications (e.g. having ChatGPT call your code or APIs). (See OpenAI's Function Calling guide for details.)

Common Pitfalls & How to Avoid Them

- **Vague or Overly Broad Prompts:** Beginners often write very general prompts and expect the model to guess their intent. As one guide warns, “A vague or poorly worded prompt can leave you with something generic or off-base” ⁴. *Solution:* Be as explicit as possible. Specify exactly what format or content you want, include context clues, or restrict the domain. For example, instead of “Tell me about physics,” say “Explain Newton's second law with a simple example.”
- **Ignoring Specificity or Role:** Not assigning roles or context can yield bland answers. If you fail to indicate domain or style, the LLM may default to generic language. *Solution:* Use role prompting or examples to anchor the response. E.g., “Act as a financial analyst” or “Translate this text as if speaking to a child” provides necessary guidance ⁷.
- **Overloading with Multiple Tasks (Overprompting):** Cramming unrelated instructions into one prompt (e.g. “Write a summary, then translate it, then list questions”) confuses the model. As one analysis notes, “overstuffing instructions in a single prompt... leads to unclear or shallow output” ¹⁷. *Solution:* Break complex tasks into sequential prompts (prompt chaining). For instance, first ask for a summary, then in a follow-up prompt ask to translate that summary. This yields better results and gives you control over each step ¹⁸.
- **Not Iterating or Refining:** Assuming a single prompt should be perfect often leads to disappointment. Most high-quality outputs come from iterative refinement. *Solution:* Examine the model's answer, then ask follow-ups or tweak the prompt. For example, if the tone is too formal, follow up with “Make it more conversational.” Palantir advises using a “feedback loop” to adjust wording and improve results ¹². Treat the LLM like a collaborator that you can guide with follow-up instructions.

- **Ignoring Model Limits / Hallucination:** LLMs generate plausible-sounding answers but can “hallucinate” facts. Over-reliance without verification is dangerous. *Solution:* Use prompts that ground the model in facts or external data. For critical info, use RAG or fact-check against sources. Remind the model to be accurate (e.g. “Only use information from the text below”). If accuracy is vital, always treat LLM output as a draft that needs human review.
- **Context-Length Issues:** Exceeding the model’s context window causes truncation or forgetting earlier parts of the prompt. *Solution:* For very long inputs, summarize or chunk the text. You might pre-process a large document into key points, or use retrieval to fetch only relevant segments ¹⁴. Some systems support “long-context” models, but generally keep prompts focused and within limits.
- **Prompt Injection / Security Risks:** Untrusted input could manipulate system prompts. For example, an attacker might craft a user query that causes the model to ignore earlier instructions. *Solution:* Sanitize user inputs and keep critical instructions in “system messages” that the user can’t override. Employ detection or filtering tools (see OWASP GenAI LLM01 for mitigation strategies) ¹⁹. Restrict models’ permissions and continually update safety checks.

Real-World Examples by Domain

- **Code Generation:** LLMs excel at translating instructions into code. For example, telling GPT-4 “Write a Python function to reverse a linked list” can produce working code. Professional developers use prompt engineering in tools like GitHub Copilot, which maps user intent to code suggestions ²⁰. Effective coding prompts specify language, function names, and requirements (e.g. “in-place, iterative solution”). Using technical context (like referring to existing classes) further improves output. In practice, one might chain prompts (“First, write this helper function. Next, write the main function.”) to build complex code step by step ¹⁸.
- **Data Analysis:** LLMs can help generate analysis plans or even code for data processing. A prompt like “Analyze the following sales dataset for Q1 and output key insights” might yield Python code or a summary. Role prompting helps (e.g. “You are a data scientist summarizing these results for the CFO”). Often teams provide schema or column descriptions in the prompt. Combining with tools (like using LLMs to generate SQL queries or Pandas scripts) is common. For instance, asking “Given this JSON dataset, write Python code to compute average order value” yields a helpful code snippet. Chaining and examples (“Show me how to compute X, then Y”) aids clarity.
- **Summarization:** LLMs are widely used to condense documents. A prompt such as “Summarize the key findings of the research paper below in plain English” leverages their abstractive summarization skill. In domain-specific contexts (medical, legal), prompts often include instructions to simplify jargon. For example, Microsoft’s team showed that adding context and roles (“You are a medical researcher writing a plain-language summary” ²¹) dramatically improves results on technical papers. One best practice is to give explicit instructions about audience and format (e.g. “Write a 200-word summary for a general audience, avoiding technical terms” ²¹).
- **Customer Support:** Chatbots use prompt engineering to handle inquiries. A typical approach is to prompt the agent with conversation context plus role: e.g. “You are an empathetic support agent. A customer says: ‘I can’t log in.’ Respond politely and helpfully.” Including the user’s message and any relevant company policy (as context) helps generate an appropriate reply. Embedding customer data

or FAQ snippets via RAG can ensure the bot has the right info. Iterating on tone (“Use a friendly tone, no more than 50 words”) is common practice.

- **Education:** LLMs serve as tutors or content creators in education. For example, a prompt might say “You are a math tutor. Explain the concept of derivatives to a high-school student” to get an age-appropriate explanation. Role prompts ensure the answer suits the learner’s level. Prompts can request step-by-step guidance (“Solve this algebra problem and show each step”), or creation of quizzes (“Generate 5 multiple-choice questions on World War II”). Embedding [51], note how the model is positioned as a teacher guiding a student. Studies show that specifying expertise (like “act as a physics professor”) leads to more accurate and tailored educational answers ⁷.

Tools & Frameworks

- **LangChain (Python/JS):** A popular open-source framework for building LLM applications. LangChain provides abstractions for prompt templates, chains, and integrations (e.g. with databases or APIs) ²². It simplifies common patterns like RAG (document loaders + vector stores + LLM) and agentic behavior (ReAct agents). Using LangChain’s prompt template classes helps manage complexity and reuse prompts. (See IBM’s overview: “LangChain... simplifies the process of building LLM-driven applications like chatbots and AI agents” ²².)
- **OpenAI Function Calling (Tool Use):** As mentioned, OpenAI’s API now supports “function calling,” which lets the model return structured JSON to invoke predefined functions. This is supported in the OpenAI SDK and simplifies integrating LLMs with external tools. (Other vendor LLMs have similar features.) Prompt engineers often use function calling to enforce structure or limit output to known schema, e.g. for databases or calculators.
- **Prompt Injection Protection:** Security tools and best practices are emerging. OWASP’s GenAI project has cataloged prompt injection vulnerabilities ¹⁹. Some LLM development kits (like LangChain) allow “prompt templates” that separate user content from system instructions, reducing injection risk. There are also detection libraries (like mitigators that scan for malicious patterns) and services offering LLM safety modules. In enterprise settings, “semantic kernel” (Microsoft) and similar frameworks help lock down prompts.
- **ReAct (Reason+Act) Agents:** ReAct is a prompting framework where the model interleaves reasoning steps with external actions ²³. For instance, a medical question could prompt: “*Question: ...*\n*Thought: I should look up this symptom.*\n*Act: search('symptom guidelines')*\n*Observation: [result]...*\n*Thought: The answer is...*” This pattern, popularized by Yao et al. (ICLR 2023), lets LLMs consult tools (APIs, search) mid-prompt ²³. Libraries like LangChain support building ReAct-style agents. Using ReAct is a powerful prompt technique when you need the model to fetch data or perform dynamic reasoning.
- **Semantic Search / Vector Stores:** Often used in retrieval pipelines (RAG). Tools like Pinecone, Weaviate, or open-source vector DBs index embeddings of documents. LangChain and others provide abstractions (Embedding + VectorStore). When prompt length is tight, you retrieve relevant passages by semantic search and include them in the prompt. This ensures only pertinent context is fed to the model.

Learning Resources

- **Prompt Engineering Guides & Repositories:** The [Prompt Engineering Guide](#) (GitHub) is a community-curated collection of papers, tutorials, and tools. As Lilian Weng notes, this repo “contains a pretty comprehensive collection of education materials on prompt engineering” ²⁴. Similarly, the *OpenAI Cookbook* (GitHub) provides example notebooks and prompt patterns.
- **Tutorial Sites & Blogs:** *LearnPrompting.org* is a free, open curriculum (with courses) on prompting basics. *PromptingGuide.ai* offers in-depth articles on specific techniques (zero-shot, CoT, etc.). Lilian Weng’s blog (“Lil’Log”) has a detailed *Prompt Engineering* post (updated 2023) which surveys many methods. Microsoft’s DevBlogs and others (e.g. Google Cloud’s Vertex AI docs) also publish practical tips. Example: Microsoft’s “Prompt Engineering for Complex Summarization” shows step-by-step prompt refinement ²¹.
- **Courses & Tutorials:** DeepLearning.AI offers a *ChatGPT Prompt Engineering* course (free audit) covering fundamentals and use cases. Kaggle and Coursera have short Guided Projects on prompt engineering. YouTube creators (DataCamp, edureka, etc.) provide walkthroughs (e.g. designing prompts for code or analysis). Many are free or low-cost.
- **Community Repos & Lists:** Search GitHub for “awesome-prompt-engineering” or “prompting resources.” Examples include [prompts-lab/Awesome-Prompt-Engineering](#). OpenAI’s own documentation (Best Practices guide) is also informative. Keeping an eye on conferences like ACL/EMNLP, NeurIPS or ArXiv (for papers on CoT, RAG, etc.) helps track the latest techniques.
- **Tools Documentation:** Read the docs of frameworks and APIs. LangChain’s docs (Python/JS) include practical tutorials (e.g. building RAG). OpenAI’s API docs cover function calling and system message usage. OWASP GenAI has a free guide on prompt injection (LLM01:2025).

Expert Tips

- **Iterate and Test:** Treat prompts like software: test them and refine. Don’t expect perfection on the first try. Use A/B testing or evaluations (automated or with humans) to compare prompt variants. Palantir emphasizes adjusting prompts based on output quality ¹².
- **Modular Prompts:** Break tasks into steps (prompt chaining) rather than one monolith. This not only avoids overloading the model ¹⁸ but also lets you control and debug each part. For complex outputs, first generate an outline or parse data, then refine it in later prompts.
- **Document Good Prompts:** As you experiment, save and catalog effective prompt patterns. Create templates with placeholders for variables. Version control (e.g. in Git) helps track changes. This practice scales well in teams and ensures reusability.
- **Stay Informed:** The field is evolving fast. Follow prompt engineering blogs (e.g. Lilian Weng, AI newsletters) and read new papers (e.g. on CoT, RAG, few-shot). Join communities like the Learn Prompting Discord or forums. New model features (like retrieval plugins or multimodal prompts) appear regularly.

- **Know Your Model:** Different LLMs have different strengths. For example, some excel at reasoning (GPT-4), others at factual recall (Claude). Tailor your prompt length, style, and expectations to the model you use. Always check the model's documented capabilities and limitations.
- **Safety and Ethics:** Build in guardrails. Use content filters or safety checks if deploying prompts in user-facing apps. Prompt injection should be top-of-mind – keep critical instructions in system prompts and validate model outputs, especially in high-stakes domains. As OWASP notes, well-designed system prompts and ongoing monitoring are needed to mitigate risks ¹⁹.

Summary: Mastering prompt engineering means learning by doing and staying current. Use clear, detailed instructions; iterate rapidly; leverage available tools (LangChain, API features, evaluation frameworks); and follow the latest community best practices ⁵ ¹². Over time, you'll build an internal library of prompt strategies and a sharper intuition for how best to interact with LLMs.

Sources: See references for all citations ¹ ⁷ ⁸ ²⁴ ¹² ¹⁹ and others. The image captions are credited automatically by the UI.

¹ ⁵ ¹¹ Prompt engineering best practices for ChatGPT | OpenAI Help Center

<https://help.openai.com/en/articles/10032626-prompt-engineering-best-practices-for-chatgpt>

² ⁹ ¹⁵ ²⁴ Prompt Engineering | Lil'Log

<https://lilianweng.github.io/posts/2023-03-15-prompt-engineering/>

³ ⁶ ¹⁰ ¹² ¹³ ¹⁴ Best practices for LLM prompt engineering • Palantir

<https://www.palantir.com/docs/foundry/aip/best-practices-prompt-engineering>

⁴ ¹⁷ ¹⁸ 5 Common Prompt Engineering Mistakes Beginners Make

<https://www.mygreatlearning.com/blog/prompt-engineering-beginners-mistakes/>

⁷ Assigning Roles to Chatbots

<https://learnprompting.org/docs/basics/roles?srsId=AfmBOor4cBcyXXJgTaaH6u7ogkninssKz5kA7dfO-ejZiv1JlpiAdbz1>

⁸ [2201.11903] Chain-of-Thought Prompting Elicits Reasoning in Large Language Models

<https://arxiv.org/abs/2201.11903>

¹⁶ RAG vs fine-tuning vs. prompt engineering | IBM

<https://www.ibm.com/think/topics/rag-vs-fine-tuning-vs-prompt-engineering>

¹⁹ LLM01:2025 Prompt Injection - OWASP Gen AI Security Project

<https://genai.owasp.org/llmrisk/llm01-prompt-injection/>

²⁰ A developer's guide to prompt engineering and LLMs - The GitHub Blog

<https://github.blog/ai-and-ml/generative-ai/prompt-engineering-guide-generative-ai-llms/>

²¹ Large Language Model Prompt Engineering for Complex Summarization - ISE Developer Blog

<https://devblogs.microsoft.com/ise/gpt-summary-prompt-engineering/>

²² What Is LangChain? | IBM

<https://www.ibm.com/think/topics/langchain>

²³ ReAct Prompting | Prompt Engineering Guide

<https://www.promptingguide.ai/techniques/react>