

Prompt Engineering	2	
1. Put instructions at the beginning of the prompt		4
2. Be Clear and Specific	5	
3. Role Prompting	5	
4. Use of Delimiters	6	
5. Control Output Format	7	
6. Focus on positive instructions	8	
7. Use Contextual Information	8	
8. Experiment with Prompt Formats	9	
9. Control the Output Length	9	
10. Asking model to explain its reasoning	10	
11. Iterate and Refine	10	
For complex prompts	11	
Prompt Engineering Types	13	
Zero Shot Prompting	13	
Few Shot Prompting	14	
Many-shot Prompting	15	
In-Context Learning	15	
Metadata Prompting	16	
Chain of Thought (CoT) prompting	17	
Prompting for multi-step processes	19	
Prompt chaining	20	
Stepwise prompting	20	
Illustration of Stepwise prompting and Prompt Chaining: Generate document outline		21
Prompt Chaining:	21	
Stepwise prompting:	22	
Missing paragraphs or lines	23	
Stepwise prompting vs Prompt Chaining	23	
Recommendation for choosing between the two:		24
Format of Generated Responses	24	
Markdown & Table Markdown	26	
Prompt Engineering Pipeline and Evaluation	28	
Pydantic Library	29	

Prompt Engineering

Prompt engineering involves crafting inputs (prompts) that guide AI models, like ChatGPT, to generate desired outputs or perform specific tasks efficiently.

Prompts for AI models have four key components:

Instruction/Intent:

It specifies what needs to be done, conveying the intent of the task to the model. This guides the model toward a specific action or response type. Clear instructions help the model focus on the task at hand. For example, instead of a vague instruction like "Create a blog post about AI," a more precise instruction would be "Write a 1000-word blog post about how AI is transforming industries. Structure the post to cover: (1) AI's impact on healthcare, focusing on diagnostic accuracy and drug discovery, (2) the revolution in manufacturing through predictive maintenance and automated quality control, and (3) AI's role in personalized education, with examples of adaptive learning platforms. Include relevant statistics and real-world case studies for each section." Such detailed instructions provide explicit guidance about content, structure, length, and specific elements to include. This precision not only reduces ambiguity but also enables the model to produce more focused, comprehensive, and well-organized responses that align exactly with the intended purpose.

Context/Data:

It provides the background information or data on which the model bases its response. It serves as the foundation for reasoning and answering questions. For instance, if the instruction is "Provide a concise summary of this scientific article highlighting the key findings and methodology in no more than 200 words," the context would be the full scientific article that the model needs to analyze and condense. This relationship between instruction and context demonstrates how the instruction defines the task's parameters while the context provides the raw material to work with. Context helps the model understand the scenario or subject matter, ensuring its responses are informed and relevant. In cases where external information is required, techniques like Retrieval-Augmented Generation (RAG) can be employed. RAG retrieves relevant context from a knowledge base or external sources based on the query, allowing the model to incorporate pertinent information dynamically. However, it's important to note that context isn't always necessary - in many cases, a well-crafted instruction alone can be sufficient. For example, an instruction like "Write a Python function to calculate the Fibonacci sequence up to n terms" doesn't require additional context since the task is self-contained and draws on the model's pre-trained knowledge of programming concepts.

Output format:

The output format specification defines the structure and presentation of the model's response. For example, instead of getting a freeform text response, you might specify "Return the competitor analysis in JSON format with the following structure: { 'market_position': { 'strengths': [...], 'weaknesses': [...] }, 'revenue_streams': { 'primary': [...], 'secondary': [...], 'growth_rate': number }, 'key_differentiators': [...] }." This level of format specification ensures not only consistent structure but also enables direct integration with dashboards or automated reporting systems. Similarly, you might request "Generate a Mermaid flowchart to visualize the user authentication process" or "Provide the

technical documentation in Markdown format following GitHub-flavored Markdown conventions." This specification ensures the output is not just accurate but also readily usable for its intended purpose, whether that's system integration (JSON/XML), documentation (Markdown), or visual representation (Mermaid/SVG).

Reasoning assistance / Guidance:

Reasoning assistance influences how the model processes the instruction and context to arrive at a response. Without this guidance, the model relies solely on its pre-trained knowledge. However, for tasks requiring domain expertise, reasoning guidance provides additional structure through two main approaches. The first is rule-based guidance, where expert-defined rules constrain the model's reasoning process and help model benefit from experts' tacit knowledge—the intuitive understanding and professional insights accumulated through years of practical experience. The second approach applies when decision rules are not explicitly clear - here, we can use traditional machine learning models like decision trees to analyze historical data and uncover underlying patterns and decision criteria. For example, in medical diagnosis, while some conditions have clear diagnostic protocols, others might have subtle patterns that can be discovered by analyzing thousands of past cases using decision trees or similar interpretable models. These learned patterns can then be incorporated into the guidance to enhance the model's reasoning process.

These four components—instruction, context, reasoning guidance, and output format—work together to make generative AI more effective, adaptable, and aligned with human expectations.

Instruction

Create a blog post about AI referring to the following discussion

Reasoning Guidance

Use the following steps:

- 1) Analyze the discussion provided
- 2) Create outline of the article
- 3) Write the full article
- 4) Write metadata

Output Format

See that blog post has appropriate headings. And also provide metadata for the post.

Refer to the following example:

```
<title>How to Start a Successful Blog in 2023</title>
<meta name="description" content="Learn step-by-step strategies to create a profitable
blog, including keyword research, content creation, and SEO optimization." />
<meta name="keywords" content="blogging, SEO, content marketing, website creation, online
business" />
```

Context/Data

Discussion:

AI has been a game-changer in my work and learning. Tools like Cursor, Midjourney, Vercel, Anthropic's Claude, OpenAI can save time and provide clever solutions. But here's the thing: I've also seen the dark side, and it's not the AI itself—it's how we use it. Enter the "AI zombies." These are folks (and maybe we've all been guilty at some point) who lean so heavily on AI that they lose touch with the critical skills that make us, well, human: problem-solving, creativity, and the ability to think deeply.

Here's a moment that really hit me:

A study by Uplevel revealed something startling. Developers who relied on tools like Copilot introduced 41% more bugs into their code. And the kicker—these same developers weren't more productive than their manual-working counterparts. I've seen this play out firsthand too: students or new hires blindly trusting AI-generated answers without a second thought. No debugging, no validating, no asking why. It's like handing over the wheel to a autopilot without learning how to drive.

Next, we will discuss details of creating prompts effectively, offering tips to enhance interactions with AI models. It explains how to write prompts that result in better, more precise, and interesting AI responses.

1. Put instructions at the beginning of the prompt

Less effective ✗:

```
{text input here}
```

Summarize the text below as a bullet point list of the most important points.

Better :

```
Summarize the text below as a bullet point list of the most important points.
```

```
{text input here}
```

2. Be Clear and Specific

The "Be Clear and Specific" guideline underlines the need for precise language when engaging with language models to ensure clear understanding and relevant responses.

Key Points:

- Clarity: Use simple language to avoid misinterpretation.
- Specificity: Provide detailed prompts to direct the AI's focus.
- Avoid Ambiguity: Use unambiguous language to prevent undesired responses.

Example 1:

A business analyst seeking to improve customer service with AI should use specific prompts like, "How does AI improve customer service in retail?" rather than vague ones like, "Tell me about AI in business." This specificity yields targeted insights.

Example 2:

A medical researcher asking about AI in cancer diagnostics should be specific: "Summarize advancements in AI diagnostic tools for cancer," to get relevant and current information.

In summary, clarity and specificity streamline AI interactions, ensuring responses are both useful and aligned with user needs, as shown in the customer service enhancement and medical research examples.

3. Role Prompting

Role prompting involves assigning a specific role or persona to the AI, which comes with a predefined set of knowledge and skills. This technique actively molds the AI's "identity" throughout the interaction. The AI takes on the communication style, level of expertise, and viewpoint associated with the assigned role. Role prompting essentially gives the AI a specific character to embody, including expected expertise and communication mannerisms. It's akin to instructing the AI to assume the persona of a particular individual with a designated background and knowledge.

Example:

If the AI is told, "You are a DevOps engineer," in a software development company, it will answer questions as if it were an actual DevOps engineer, using technical language and concepts specific to software development and IT operations. This goes beyond general thematic guidance to provide specific, expert-level responses based on the assigned role.

4. Use of Delimiters

You may use delimiters such as "" or ``` to separate the instruction and context.

Less effective ✗:

```
Summarize the text below as a bullet point list of the most important points.
```

```
{text input here}
```

Better ✓:

```
Summarize the text below as a bullet point list of the most important points.
```

```
Text: """
```

```
{text input here}
```

```
"""
```

Some recommend using ### or - - - in between two examples in your prompt. Delimiters in prompt engineering help improve Large Language Models (LLMs) by guiding their attention to specific parts of the input. Delimiters, such as triple quotes, backticks, dashes, angle brackets, and XML tags, separate important sections of the prompt, making it easier for the model to prioritize essential information.

Triple quotes: """

Triple backticks: ```

Triple dashes: - - -

Angle brackets: <>

XML tags: <tag> </tag>

This approach not only enhances the accuracy and relevance of the model's responses but also safeguards against prompt injection, where the model might unintentionally execute misleading instructions.

For example, imagine you need a concise summary of a market analysis from a lengthy report in a business scenario. Without delimiters, requesting a summary might result in a broad overview, losing focus. By enclosing the market analysis section with triple backticks in your prompt, you signal the LLM to concentrate on that specific section, leading to a targeted summary. Using delimiters ensures

more precise and helpful outputs from the model by clarifying your instructions and directing its attention, thus improving its effectiveness.

Prompt:

Please summarize the following section of the report:

'<market analysis section of the report>'

The prompt directs the AI to summarize the specified "market analysis section" of the report, ensuring a focused and useful summary by making instructions explicit and guiding attention to critical content, thus enhancing performance.

5. Control Output Format

Many times you would need the response from a foundational model such as ChatGPT in a specific format such as JSON, or Python lists to ensure seamless integration with other systems and applications. This approach facilitates easy parsing and use of the data in automated processes, eliminating the need for manual formatting.

Imagine, you are developing a travel app that suggests restaurants based on dietary preferences. Based on user's responses and preferences, let's say you created a prompt that generates restaurant recommendations. However, without any specifics about the format, LLM is going to return recommendations as unstructured text, which would be difficult for you to extract. You can completely avoid this extra work by specifying the format of the response in the prompt like the following:

Prompt: "Provide a list of three best restaurants near the Eiffel Tower in Paris that offer vegetarian food options. Your response should be in JSON format, [{restaurant:<restaurant name>, price: <expected price per head>}]."

Example response structure (simplified for clarity):

```
[  
  {"restaurant": "Vegetarian Delight Bistro", price: 25},  
  {"restaurant": "Green Garden Café", price: 43},  
  {"restaurant": "Eco Eats", price: 18}  
]
```

Occasionally, the foundational model might not deliver the output in the specified format you've requested. You can fix that by using few-shot prompting, ie providing sample output along in the prompt.

Prompt: Provide a list of three best restaurants near the Eiffel Tower in Paris that offer vegetarian food options. Your response should be in JSON string format, [{restaurant:<restaurant name>, price: <expected price per head>}].

Sample response:

```
[  
  {"restaurant": "restaurant 1", price: 12},  
  {"restaurant": "restaurant 2", price: 23}
```

]

6. Focus on positive instructions

The guideline of focusing on positive instructions when crafting prompts for ChatGPT suggests a direct and affirmative way of communicating what you want the AI to do. This method is preferred because it clearly states the expected outcome, reducing the chance of misunderstanding and making it easier for the AI to follow the instructions accurately.

Example 1:

Say you want to write a product description for a tea pot based on the things customers are writing in their review. You would want to make sure the description doesn't discuss anything negative about the tea pot.

In such a scenario, instead of negative prompt instruction such as:

"Write a product description for the [model] tea pot based on the following customer reviews. Do not write description based on negative reviews."

Use a Positive Prompt Instruction such as:

"Write a product description for the [model] tea pot based on the following customer reviews. Write description based on positive reviews only."

Example 2:

Instead of the following:

"Write a title for the following article with no emotional tone"

Use:

"Write a title for the following article with formal tone"

7. Use Contextual Information

The guideline "Use Contextual Information" is about making ChatGPT prompts more effective by adding details specific to the user's situation, such as where they are, what they like, and interactions they've had before. This makes the responses from ChatGPT more tailored and directly relevant to the user.

Example:

Imagine a travel application that incorporates ChatGPT for restaurant suggestions. By integrating prompts filled with contextual details, such as the user's exact location and their dietary preferences, the app is capable of offering custom-tailored restaurant options. This approach to personalization significantly improves the user experience, simplifying the process of finding suitable dining establishments that cater to the user's unique needs. Utilizing context-rich prompts, the app delivers more relevant and practical recommendations, demonstrating the effectiveness of adding specific, user-related information into AI-powered tools to enhance their usefulness and accuracy.

Prompt:

"Given a user's current location and their preference for vegetarian dining options, generate a list of Italian restaurants within a 5-mile radius of their position. Include details such as the restaurant's name, address, a brief description of its vegetarian offerings, and any notable features (e.g., outdoor seating, reservation requirements). Ensure the recommendations are tailored to the user's specific request, emphasizing the quality and variety of vegetarian dishes available."

8. Experiment with Prompt Formats

The "Experiment with Prompt Formats" guideline emphasizes the need to vary prompt structures when using AI, as this significantly affects the AI's responses. Experimenting with questions, statements, or instructions helps tailor AI responses to better meet expectations.

Consider a content writer in the sustainable fashion industry aiming to produce content on environmental conservation. They can use:

1. Question Prompt: Asking, "How does sustainable fashion contribute to environmental conservation?" prompts a broad, exploratory response, ideal for understanding sustainable fashion's overall impact.
2. Statement Prompt: "Describe the role of sustainable fashion in reducing the fashion industry's environmental impact" yields a focused explanation on specific impacts, suitable for detailed articles.
3. Instruction Prompt: "Provide three examples of sustainable fashion brands and their eco-friendly practices" directs the AI to list specific brands, useful for creating content with real-world examples.

Experimenting with these formats enables the writer to gather diverse information, enriching content and tailoring it for various formats like blog posts, in-depth articles, or listicles on eco-friendly brands.

In essence, varying prompt formats enhances the use of AI tools, making them more aligned with specific content creation goals, especially in fields aiming to produce engaging, informative content for a target audience.

9. Control the Output Length

The "Control Output Length" guideline allows users to manage the length of ChatGPT responses through two methods: setting a character limit or requesting a specific length (e.g., number of sentences). For instance, a character limit of 280 characters suits brief formats like tweets, while asking for a response "in three sentences" ensures conciseness without sacrificing completeness.

Illustrating a business example, a social media manager promoting a new smartphone on Twitter, which has a 280-character limit, might use the prompt: "Describe the key features of our new smartphone in a tweet-length message." This ensures the message is concise, focuses on essential

features like camera quality and battery life, and fits within a single tweet. This guideline aids in producing directly usable content, saving time and effort in editing, and effectively conveying the product's value within platform constraints.

10.

Asking model to explain its reasoning

Asking an AI model to generate explanations for its labels or recommendations can significantly enhance output quality by promoting deeper reasoning and analysis. This approach, closely related to chain-of-thought prompting, encourages the model to articulate its decision-making process, which can reveal and potentially correct flaws in its reasoning. By requiring explanations, the model is pushed to engage in more thorough contextual understanding and align its thinking more closely with human-like reasoning patterns. This process can help mitigate biases, improve transparency, and ultimately lead to more thoughtful, well-justified outputs. Additionally, the act of explaining can reinforce the model's grasp of concepts and relationships, potentially improving its performance over time. This technique not only enhances the model's ability to handle complex tasks but also provides valuable insights into its decision-making process, fostering greater trust and understanding between AI systems and their users.

11. Iterate and Refine

The "Iterate and Refine" guideline in prompt engineering highlights the necessity for continuous testing, evaluation, and enhancement of prompts used with AI models, like ChatGPT, to optimize response efficiency and accuracy. This iterative process involves experimenting with various prompts, analyzing AI responses, and refining prompts based on performance to improve response quality and relevance gradually.

Acknowledging the trial and error involved is essential, as crafting the perfect prompt often requires multiple attempts due to the complexities of human language and AI interpretation. Initial attempts may not fully convey the needed context or specificity, necessitating prompt adjustments.

- Ask for LLM's understanding of the prompt

Start by ensuring the AI comprehends your prompt correctly. This step involves not just asking for understanding, but also an iterative refinement process. Here's the expanded process:

a. Initial Query:

Use this specific prompt to get the AI's initial understanding:

Provide your understanding of the following prompt for an AI tool:

b. Analyze the Response:

Carefully review the AI's explanation of your prompt. Look for any misinterpretations, gaps in understanding, or areas where the AI's interpretation doesn't align with your intent.

c. Iterative Refinement: Use the edit option to change your original prompt. Update the prompt to incorporate better wordings or explanations you see in the AI's output. When you save the updated prompt, the AI will give you another explanation. Review this new explanation carefully.

d. Decision Point: If you see the need for further minor changes, repeat the process from step c. If you're satisfied with the AI's understanding and feel no further changes are necessary, proceed to the next step in the prompt engineering process.

It may take you 2-3 iterations to fix your prompt. This iterative refinement within the first step is crucial because it allows you to:

- Gain insights into how the AI interprets your language
- Incrementally improve your prompt based on the AI's feedback
- Ensure a solid foundation of mutual understanding before moving on to more complex refinements

An alternative to manual intervention is to let LLM handle the rewrite. Once you've confirmed that the AI's understanding is good and it hasn't misconstrued or deviated much from your intent, ask it to improve the prompt:

Rewrite the prompt to make it better

Or

Evaluate the structure of the following content, focusing on improving its organization and presentation. Avoid adding or suggesting new information--your task is to reframe the existing content for better clarity and flow.

This collaborative approach can lead to unexpected insights and refinements.

-----NOT part of the Exam-----

For complex prompts

The first draft of your prompt can be approached like any writing task - simply jot down your initial thoughts without concern for structure or order. This approach is particularly beneficial for complex prompts covering extensive ground and specifications. Once you've captured your ideas, you can refine the prompt using a follow-up instruction. This subsequent prompt should request a rewrite that incorporates specific considerations while preserving all original information. These considerations include improving structure and clarity by reorganizing dense paragraphs into clear headings and bullet points, optimizing the order of operations for logical flow, separating content requirements from output format specifications, ensuring consistency in formatting and style throughout, and eliminating redundancy by consolidating repetitive instructions. This process will result in a more refined, clear, and concise prompt that maintains the integrity of your original ideas

while enhancing its overall effectiveness.

Rewrite the prompt and incorporate the following considerations. It is imperative that you do not leave out any information in the original prompt.

- Structure and clarity: If the current prompt is dense and combines multiple instructions in long paragraphs, restructure it with clearer headings and bullet points to make it easier to follow and implement.
- Order of operations: If the steps aren't presented in the most logical order, fix the order.
- Separation of concerns: If the prompt mixes content requirements with output format specifications, separating these could improve clarity.
- Consistency: If the formatting and style of instructions vary throughout the prompt, make its style more consistent to improve readability.
- Redundancy: If there's some repetition, particularly in the explanation requirements, consolidating these instructions would make the prompt more concise.

- Addressing potential uncertainties

Next, address any subjectivity or unclear elements in your prompt that could lead to unreliable results, especially when the context might differ from your test cases.

- Identify potential uncertainties

You can begin with using LLM to help you identify uncertainties in your instructions by using the prompt given below. This question helps you pinpoint areas where your prompt might be open to interpretation or lacking specificity.

Is there any subjectivity in the prompt or something unclear for an AI tool

- Asking LLM to guess answers for uncertainties

Instead of manually addressing how to make instructions more specific for identified uncertainties, you can ask the LLM to make educated guesses about potential answers or solutions. This approach capitalizes on the model's advanced capabilities, potentially saving you time and effort. By using the prompt given below you're essentially outsourcing part of the problem-solving process to the AI. This not only helps in generating potential solutions but also provides insights into how the model might interpret and respond to ambiguities in your prompt, further informing your refinement process.

Make your best guess and try to answer subjectivities you identified in the last response

LLM can overdo and list frivolous points at times, besides mostly great feedback on uncertainties. You would want to filter good points from the rest.

- Ask LLM to rewrite prompt to address uncertainties

Based on the insights gained, you can ask LLM to rewrite your prompt to address the identified issues. Remember that the AI might overdo it and list some frivolous points alongside mostly great feedback. You may need to mention specific points that you want to be incorporated leaving the rest. By selectively incorporating points, you prevent the prompt from becoming overly complex or veering off-track due to the AI's tendency to sometimes over-elaborate.

Rewrite the prompt to address the following points:

- point 1
- point 2

While these guidelines provide a solid foundation, don't hesitate to experiment with different phrasings, structures, and approaches. Each use case may require unique tweaks to achieve optimal results. By combining thoughtful design with systematic testing and refinement, you can create highly effective prompt templates that maximize the capabilities of LLMs in your workflow.

Prompt Engineering Types

Zero Shot Prompting

Zero-shot prompting is a technique used with Generative Pre-trained Language Models (LLMs) like GPT (Generative Pre-trained Transformer) that enables the model to undertake tasks it hasn't been explicitly trained on. It involves presenting a task to a language model without any task-specific examples or training. The model is expected to understand and execute the task based solely on its pre-existing knowledge and the general instructions provided in the prompt. We communicate with the model using a prompt that explains what we want to achieve. The model uses its pre-trained knowledge, acquired from a vast amount of text data, to infer the best way to complete the task.

This capability is pivotal for several reasons:

- Versatility and Adaptability: It allows models to handle a wide range of tasks without the need for fine-tuning or retraining, making them highly versatile and adaptable to new challenges. Whether it's sentiment analysis, summarization, or question-answering, the model adapts to the prompts provided.
- Cost Efficiency: Reducing the necessity for large, annotated datasets for every new task saves significant resources in data collection and annotation.
- Generalization: Demonstrates the model's ability to generalize from its training data to new, unseen tasks, highlighting its understanding of language and concepts.

Example of Zero-Shot Prompting

Let's consider the task of sentiment classification. Here's how you would set up your prompt:

Task: Sentiment classification

Classes: Positive, neutral, negative

Text: "That shot selection was awesome."

Prompt: "Classify the given text into one of the following sentiment categories: positive, neutral, negative."

The model's response would likely be "positive" because it has learned from its training data that the word "awesome" is associated with positive sentiment.

Few Shot Prompting

Few-shot prompting is a technique used to guide large language models (LLMs), such as ChatGPT and Llama, to perform specific tasks or understand particular contexts using only a small number of examples. In few-shot prompting, you provide the model with a few carefully selected examples (typically between 2 and 10) that demonstrate both the input and the desired output of the task. These examples help the model infer the pattern or context of the task, which it then attempts to generalize to new, unseen inputs.

It's important to note that the model does not update its internal weights during few-shot prompting. The model temporarily "learns" or infers patterns from the provided examples but discards this information once the interaction is over.

Example 1:

Input: "Do you have the latest model of the XYZ smartphone in stock?"

Response: "Thank you for your inquiry. Yes, we have the latest XYZ smartphone model available. Would you like to place an order?"

Example 2:

Input: "Is the ABC laptop available in your store?"

Response: "Thank you for reaching out. The ABC laptop is currently out of stock, but we expect new shipments to arrive next month. Can we notify you when it's available?"

Your task:

Input: "Can you tell me if you have the DEF headphones in stock?"

Response:

In this scenario, the model is provided with two examples of customer inquiries regarding product availability, along with the corresponding email responses. In the first example, the product is in stock, and the response includes an offer to place an order. In the second example, the product is out of stock, and the response offers to notify the customer when it becomes available.

When the model is tasked with generating a response to a new inquiry about DEF headphones, it applies the pattern observed in the previous examples to craft an appropriate reply. This might involve confirming the product's availability and suggesting next steps if it's in stock, or explaining that the product is out of stock and offering alternatives or a notification service.

This approach enables the model to understand the context of customer service in a business setting and to generate responses that are both relevant and considerate of the customer's needs.

Exemplars (Examples)

Exemplars are specific instances or examples that demonstrate how a task should be performed, helping to train or guide machine learning models, especially in few-shot learning scenarios. Here's how few-shot prompting can be approached using exemplars for a business-related task, such as drafting email responses to customer inquiries about product availability, while paying attention to avoiding common pitfalls:

- Ensure Exemplar Consistency: All exemplars should follow a consistent format and structure. This consistency helps the model to understand the task and apply its learning to new inputs effectively.
- Select Relevant Exemplars: Choose exemplars directly related to the task at hand. Irrelevant exemplars can confuse the model, leading to inaccurate outputs.
- Diversify Your Exemplars: To give the model a broad understanding of the task, include a range of exemplars that cover various scenarios and outcomes related to the task. This diversity helps the model handle different inputs more effectively.
- Keep Exemplars Simple and Clear: While it's important to capture the complexity of the task, overly complicated exemplars can confuse the model. Aim for clarity and simplicity to ensure the model can easily learn from the examples provided.
- Optimize the Number of Exemplars: Balance is key. Too few exemplars may not provide enough information for the model to understand the task, while too many can overwhelm it. Adjust the number of exemplars based on the task's complexity and the model's performance.
- Incorporate Contextual Clues in Exemplars: Providing clear instructions and relevant context within your exemplars is crucial. These clues help the model to understand the task better and generate more accurate outputs.

Many-shot Prompting

Many-shot prompting is a variant of few-shot learning where, instead of using a handful of examples (e.g., around 10), you use several hundred examples (e.g., 500-800). Models with large context windows, such as Gemma, can accommodate many examples in a single prompt. However, a significant downside of utilizing such large context windows is the increased computational cost and slower inference times. With this many examples, it may be more efficient to fine-tune the model directly, avoiding the repeated cost of processing large context lengths during every inference.

In-Context Learning

In-Context Learning refers to a large language model's ability to perform tasks by interpreting examples provided in the input prompt, without updating its internal parameters. Few-shot prompting and many-shot prompting are both forms of in-context learning. Despite the term "learning," the model doesn't actually update its weights or retain information beyond the current interaction. Instead, it temporarily infers patterns or rules from the examples in the prompt but discards this inferred knowledge once the interaction concludes.

Metadata Prompting

Metadata prompting is an approach designed to simplify and streamline the process of instructing large language models (LLMs). It applies principles of modularity and separation of concerns to prompt engineering, enhancing the effectiveness of communication with LLMs. Traditionally, prompts often combine task descriptions with explanations of various entities involved, resulting in complex and cluttered instructions.

The core principle of metadata prompting is to separate the task description from entity explanations. It encourages users to start by clearly defining the main task, using all necessary entities without worrying about explaining them. To distinguish entities within the task description, they are enclosed in backticks (`). This allows for a focused and concise task description while clearly marking which terms will be explained later.

After the task is clearly defined, each entity that requires explanation is described separately in JSON format. The entity names serve as keys, with their explanations as corresponding values. This structured approach offers several benefits:

- It creates a clear separation between the task description and entity explanations.
- It makes prompts easier to understand, modify, and maintain.
- It helps visualize connections between different parts of the task more effectively.
- It reduces clutter in the main task description.
- It introduces modularity, allowing for easier updates and reuse of entity explanations across different prompts.

By structuring prompts in this way, metadata prompting aims to create more efficient, readable, and adaptable instructions for AI models, ultimately improving the quality of AI-generated outputs and making the process of working with LLMs more user-friendly.

Taking an example, let's consider a situation where a user wants to assign custom tags to each paragraph in an extensive document. Given the limitations on the token size that an LLM can handle, the document would need partitioning into segments. Yet, for every segment, crucial context like the document's title, headings, and preceding paragraphs must be provided. Traditional prompting methods might fall short here, as LLMs could have difficulty discerning metadata from the main content. In contrast, Metadata prompting offers a more straightforward communication method.

```
Tag each of `target-paragraphs` with one of the `tags` considering `article-title`, `headings` and `preceding-paragraphs`.  
tags: """  
    tagA: definition of tag A  
    tagB: definition of tag B  
""",  
article-title: """Article title""",  
headings: """  
    h1: heading with type Heading 1  
    h2: heading with type Heading 2  
"""
```

```
preceding-paragraphs: """Provide 2 paragraphs that come before the target paragraphs to give more context"""
target-paragraphs: """Provide the paragraphs you want the task to summarize"""
```

Using impressive NLU and in-context learning abilities of LLMs, AI agents typically use text as an interface between components to plan, use external tools, evaluate, reflect, and improve without additional training.

Chain of Thought (CoT) prompting

Chain of thought prompting is a technique used to encourage language models to break down complex problems into a series of smaller, interconnected steps or thoughts, mimicking the way humans reason through problems. A language model is prompted to generate a series of short sentences that mimic the reasoning process a person might employ in solving a task. The process involves three main steps:

1. Step-by-step reasoning: Instead of directly providing the final answer, the model generates a series of intermediate reasoning steps that guide it towards the solution by breaking down the problem into smaller, more manageable parts.
2. Intermediate outputs: At each step, the model generates an intermediate output that serves as a building block for the next step in the chain of thought. These outputs can be partial solutions, relevant information, or logical connections.
3. Final output: After generating the intermediate steps, the model combines the information to produce the final answer or solution to the original prompt.

There are several approaches to prompting a model to generate intermediate reasoning steps in a chain of thought. The most common and the one used in the original paper by Wei et al. (2022) is few-shot learning. In this approach, the model is provided with a few examples of problems along with their corresponding chains of thought and final answers. The model learns from these examples and applies the same reasoning pattern to new, unseen problems, relying on its ability to generalize from a small number of examples.

In their experiments, Wei et al. (2022) provided the model with examples of problems, each demonstrating the step-by-step reasoning process. For instance:

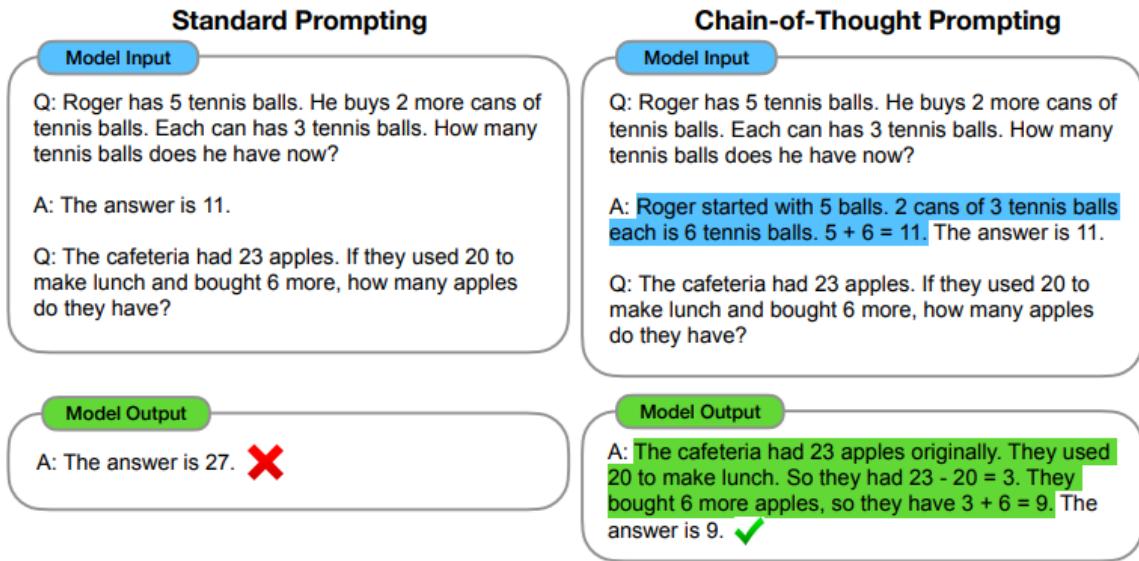


Figure 1: Chain-of-thought prompting enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. Chain-of-thought reasoning processes are highlighted.

Source: [Paper link](#)

Note: A good read for automating picking exemplars, Auto-CoT: [Paper link](#) , [Good summary article](#)

When presented with a new question, the model uses these examples as a reference to generate its own chain of thought and final answer. The authors found that this few-shot learning approach led to significant improvements in the model's performance on various reasoning tasks, including arithmetic, commonsense reasoning, and symbolic manipulation. The generated chains of thought also provided valuable insights into the model's reasoning process, making its outputs more interpretable and trustworthy.

Typical implementation:

Question 1 to n are the few shot exemplars with their respective Reasonings and Answers.

```

Question: {question 1}
Reasoning: Let's think step-by-step. {reasoning 1}
Answer: {answer 1}
...
Question: {question n}
Reasoning: Let's think step-by-step. {reasoning n}
Answer: {answer n}

Question: {question 1}
Reasoning: Let's think step-by-step.
  
```

Other approaches to prompting a model to generate intermediate reasoning steps include:

1. Zero-shot Chain of Thought: By appending the phrase "Let's think step by step", "Break down your reasoning into clear steps", or "Take a deep breath and work on this problem step-by-step" to the original prompt given to the model, it encourages the model to break down its reasoning process into a series of logical and intermediate steps rather than attempting to reach the final answer in one leap.

2. Structured prompts: Prompts that include placeholders for intermediate reasoning steps, which the model is trained to fill in along with the final answer. For instance, a prompt might be structured as follows: Question: [Original question] Step 1: [Placeholder for first reasoning step] Step 2: [Placeholder for second reasoning step] ... Step N: [Placeholder for final reasoning step] Answer: [Placeholder for final answer] The model is trained to fill in the placeholders with relevant intermediate steps and the final answer.

How is it Different from Standard Prompting?

Standard prompting might involve asking a model a direct question and receiving a direct answer, without any explanation of the steps taken to reach that answer. CoT prompting, on the other hand, explicitly asks the model to show its work, providing a step-by-step breakdown of its reasoning. This not only leads to more accurate answers in many cases but also provides an explanation that can be helpful for users to understand the model's thought process.

Business Example: Enhancing Customer Support with RAG and CoT

Consider an online retailer implementing a chatbot equipped with RAG and chain of thought prompting to handle customer inquiries. A customer asks a complicated question about a product's features, compatibility with other devices, and return policy.

- Logical Processing: Through chain of thought prompting, the chatbot first breaks down the query into sub-questions: What are the product's key features? Which devices are compatible? What is the return policy?
- Retrieval: For each sub-question, the chatbot sequentially processes the information, starting with product features, moving to compatibility, and finally addressing the return policy. At each step, it synthesizes information from the retrieved documents and previous reasoning steps.
- Final Response: The chatbot compiles its findings into a comprehensive response that clearly explains the product's features, compatibility with specific devices, and return policy, offering a detailed and helpful answer to the customer's inquiry.

This example illustrates how chain of thought prompting in RAG transforms the way LLMs handle complex queries, enabling them to provide more accurate, detailed, and contextually relevant responses. By mimicking human-like reasoning and adaptability, this approach significantly enhances the capabilities of AI in business applications, particularly in areas requiring deep understanding and nuanced responses.

Prompting for multi-step processes

Prompting often involves instructing models to produce responses for complex, multi-step processes. Try to generate prompt for the following scenario.

This scenario involves processing a document through three main steps:

- Summarizing each paragraph of the document in one line.
- Extracting key points from the entire document.
- Organizing the generated summary lines under their respective key points in the order they appear in the original document.

This task requires a combination of summarization, key point extraction, and organizational skills. It's a complex process that involves understanding the document's structure, content, and main themes.

The two popular ways to handle such tasks are: Prompt Chaining, and Stepwise Prompting.

Prompt chaining

Prompt chaining is an approach for refining outputs from large language models (LLMs) that involves using a series of discrete prompts to guide the model through different phases of a task. This method allows for a more structured and controlled approach to refinement, with each step having a specific focus. By breaking down complex tasks into smaller, more manageable prompts, prompt chaining provides greater flexibility and precision in directing the LLM's output. For example, in a text summarization task, one prompt might focus on extracting key points, another on organizing them coherently, and a final prompt on polishing the language.

Stepwise prompting

Stepwise prompting, on the other hand, integrates all phases of the task within a single prompt. This approach attempts to guide the LLM through the entire process in one go, challenging the model to generate a longer and more complex output based on a single set of instructions. While simpler to implement, stepwise prompting may be less effective for complex tasks that benefit from a more granular approach. For instance, when summarizing a lengthy academic paper, a single stepwise prompt might struggle to capture all the nuances of content selection, organization, and style that separate prompts in a chain could address individually.

Further, you can improve the clarity by assigning step names followed by step explanation. For example:

- 1) Analysis: Analyze the given text for key themes.
- 2) Summary: Summarize the themes identified in 'Analysis'.
- 3) Conclusion: Draw conclusions based on the 'Summary'.

The use of backticks to reference previous outputs leaves no ambiguity.

When to Use Step Names in Stepwise Prompts?

Use step names when the output of a step will be referenced later in the prompt or in subsequent prompts. Step names help in organization and add to clarity, especially for tasks with multiple interdependent steps. Additionally, it can help the model more effectively associate output variables with the specific steps in which they are extracted.

When Step Names Might Be Optional?

1. Simple, Linear Tasks: For straightforward tasks with clear progression, step names might be unnecessary.
2. Short Prompts: In brief prompts with only 2-3 steps, numbering alone might suffice.

Best Practices

- Be consistent: If you use step names, use them for all steps in the prompt.
- Keep names short and descriptive: Use clear, concise labels that indicate the purpose of each step.

Illustration of Stepwise prompting and Prompt Chaining: Generate document outline

Let's create prompts for the scenario explained earlier using both techniques:

Prompt Chaining:

For prompt chaining, we'll break down the task into three separate prompts, each focusing on a specific subtask.

Prompt 1 (Paragraph Summarization):

You are tasked with summarizing a document. For each paragraph in the given document, create a one-line summary that captures its main idea. Please provide these summary lines in a numbered list, with each number corresponding to the paragraph number in the original document.

Prompt 2 (Key Point Extraction):

Based on the entire document, identify and list the main key points. These should be the overarching themes or crucial ideas that span multiple paragraphs. Present these key points in a bulleted list.

Prompt 3 (Organization):

You will be provided with two lists: one containing one-line summaries of each paragraph, and another containing key points extracted from the document.

Go through the `Document`. Your task is to organize the summary lines from `Summary list` under their most relevant key points from `Key point list`. Maintain the original order of the summary lines within each key point. Present the result as a structured list with key points as main headings and relevant summary lines as sub-points. Output can be in the following format:

```
**<key point 1>
<summary line 1>
<summary line 2>
**<key point 2>
```

```
<summary line 3>
<summary line 4>
<summary line 5>
```

Ensure that all summary lines are included and that they maintain their original numbering order within each key point category.

Document: It is the document that needs to be converted into topic-wise `summary list`

Summary list: It is a list of one-line summaries of each paragraph in the `document`

Key point list: It is a list of key points covered in the `Document`

Stepwise prompting:

The prompt below demonstrates how we can replicate the three distinct steps previously used in prompt chaining within a single step by using stepwise prompting.

You are tasked with analyzing, summarizing, and organizing the content of a given document. Please follow these steps in order:

1. Read document: Carefully read through the entire document.
2. Summary list: Create a one-line summary for each paragraph, capturing its main idea. Number these summaries according to the paragraph they represent.
3. Key point list: Identify the main key points of the entire document. These should be overarching themes or crucial ideas that span multiple paragraphs.
4. Key point wise Summary list: Organize the `Summary list` under their most relevant key points from `Key point list`. Maintain the original order of the summary lines within each key point. Present the result as a structured list with key points as main headings and relevant summary lines as sub-points.

The output should be in the following format:

```
**<key point 1>
<summary line 1>
<summary line 2>
**<key point 2>
<summary line 3>
<summary line 4>
<summary line 5>
```

Ensure that all summary lines are included and that they maintain their original numbering order within each key point category.

Missing paragraphs or lines

Notice the last line: Ensure that all summary lines are included and that they maintain their original numbering order within each key point category.

Despite our repeated efforts to ensure the AI model analyzes every paragraph and line, it continues to overlook some. Models may consider multiple paragraphs as one paragraph, or may skip paragraphs completely. This is much common problem when processing long documents. Similarly, if you have to process each line in a paragraph, model may combine lines or skip lines.

The solution is you should consider breaking the document into smaller chunks and process each chunk separately. Further, I've seen better results if you provide numbered list of paragraphs, and ask it to generate output as the numbered list corresponding to each paragraph. This makes it easier for the AI model to keep track of paragraphs. Same goes for line processing within paragraphs.

Stepwise prompting vs Prompt Chaining

Here's a comparison of prompt chaining and stepwise prompting in table format:

Aspects	Prompt Chaining	Stepwise Prompting
Execution	Runs the LLM multiple times, with each step focusing on a specific subtask	Completes all phases within a single generation, requiring only one run of the LLM
Complexity and Control	Allows for more precise control over each phase of the task, but requires more comprehensive prompts from humans	Uses a simpler prompt containing sequential steps, but challenges the LLM to generate a longer and more complex output
Effectiveness	Generally yields better results, especially in text summarization tasks	Might produce a simulated refinement process rather than a genuine one, potentially limiting its effectiveness
Task Breakdown	Excels at breaking down complex tasks into smaller, more manageable prompts	Attempts to handle the entire task in a single, more complex prompt
Iterative Improvement	Allows for easier iteration and improvement of individual steps in the process	Less flexible for targeted improvements without modifying the entire prompt

Resource Usage	May require more computational resources due to multiple LLM runs	More efficient in terms of API calls or processing time
Learning Curve	Higher initial complexity for prompt designers, but potentially more intuitive for complex tasks	Simpler to implement initially, but may be challenging to optimize for complex tasks

Recommendation for choosing between the two:

I recommend starting with stepwise prompting, as it is a more cost-effective solution and requires less engineering effort compared to prompt chaining. However, if you notice a decline in quality or inconsistent results, switching to prompt chaining will be necessary.

Format of Generated Responses

Output formatting is a crucial factor in optimizing the performance of Large Language Models (LLMs). However, achieving the desired response format presents several significant challenges:

Verbosity and conversational tendencies: These models often exhibit a tendency towards verbosity and conversational behavior. In an attempt to mimic human interaction, they frequently offer unsolicited commentary on the provided prompt or their own responses, even when explicitly instructed not to do so. This chattiness can lead to unnecessary explanations that may contaminate the intended results, making it difficult to extract the precise information needed.

Formatting adherence issues: Many models, including prominent ones like Claude and ChatGPT, often struggle to strictly adhere to specified XML or JSON formats outlined in prompts. This difficulty stems from the potential conflict between strict formatting and the models' learned patterns of communication over vast amount of unstructured text data. This inconsistency in formatting can create challenges for users who require structured outputs for further processing or analysis. The inability to consistently produce responses in a desired format can hinder the seamless integration of LLM outputs into existing workflows or systems.

Tradeoff between format and reasoning: Research has revealed an interesting tradeoff between format adherence and reasoning capabilities. Imposing rigid format constraints on LLMs can actually

impede their ability to reason effectively¹². Studies have shown a correlation between stricter format requirements and a more pronounced degradation in performance on reasoning tasks. This effect may be due to the additional cognitive load placed on the model to maintain format compliance while also engaging in complex reasoning. The format-reasoning tradeoff presents a significant challenge for developers and users of LLMs. It requires careful consideration when designing prompts and evaluating model outputs, as one must balance the need for structured data with the desire to leverage the model's full reasoning potential. Often, this issue is mitigated by letting model focus on reasoning and not forcing model to stick to specific format. Information is later extracted in a structured format from the generated output.

There are several approaches to handle the problem of structuring and validating output from language models:

Output parsers

Output parsers are tools designed to process and structure the raw output from large language models (LLMs) into more usable formats. When we employ a smaller language model as part of this parsing process, we're leveraging a more specialized and often more efficient model to interpret and restructure the output of the larger, more general-purpose model. This approach involves feeding the LLM's natural language response into an SLM specifically trained for information extraction and structuring. The SLM analyzes the text, identifies key elements, and organizes them into a predefined format like JSON or XML. It can also perform data validation and refinement. This method offers benefits such as efficiency, specialization, consistency, and flexibility. SLMs are typically faster and can be fine-tuned for specific tasks, potentially improving accuracy. Frameworks like LangChain³ and LlamaIndex⁴ provide built-in support for this parsing technique. However, while effective, this approach adds complexity to the system. It's often possible to avoid using language models for data validation altogether by adopting strategies that are less demanding for the LLM to manage.

Simplified Output Format

Instructing an LLM to adhere to simplified output formats is an effective technique for obtaining more consistent and easily parseable output. This approach leverages the LLM's ability to follow instructions while simplifying post-processing. By providing clear guidelines on response formatting, we can create semi-structured output that's easier to process programmatically later without compromising the model's reasoning capabilities.

- Table Markdown: Asking the LLM to present information in a table format using markdown can be highly effective. By asking the LLM to present information in a table format using

¹ Tam, Z. R., Wu, C.-K., Tsai, Y.-L., Lin, C.-Y., Lee, H.-Y., & Chen, Y.-N. (2024). Let me speak freely? A study on the impact of format restrictions on performance of large language models. arXiv.

<https://arxiv.org/abs/2408.02442>

² Shorten, C., Pierse, C., Smith, T. B., Cardenas, E., Sharma, A., Trengrove, J., & van Luijt, B. (2024). StructuredRAG: JSON response formatting with large language models. arXiv. <https://arxiv.org/abs/2408.11061>

³ https://python.langchain.com/v0.1/docs/modules/model_io/output_parsers/

⁴ https://docs.llamaindex.ai/en/stable/module_guides/querying_structured_outputs/output_parser/

markdown, we create output that's both human-readable and machine-parseable. For example:

Name	Age	Occupation
John	30	Engineer
Lisa	28	Designer

This format is easy for the LLM to produce and for Python scripts to parse. Note: Throughout the course when I say Python script, understand that it can be done with any other language as well. It is that the data science community uses Python extensively so I'm following a popular choice.

- Prefixes or delimiters: This method involves using specific markers to distinguish different entity types. Examples include asterisks, dashes, question and answer markers (Q:, A:), arrows (=>), asterisks, or double dashes. This approach is particularly effective when dealing with a limited number of entity types.

These less rigid formatting approaches typically yield more consistent results without significantly compromising the model's reasoning abilities. If necessary, you can easily convert this prefixed output into more structured formats like JSON or XML by instructing the LLM to write a Python script for the conversion.

It's worth noting that the effectiveness of these methods can vary depending on the complexity of the data structure. For instance, when dealing with nested data or multiple values within a single field (such as storing multiple addresses), output parsers may still be the preferred solution.

-----NOT part of the Exam-----

Markdown & Table Markdown

Markdown is a lightweight markup language designed to be easy to read and write, while also being convertible to HTML and other formats. It uses simple, intuitive syntax for formatting text, creating lists, adding links, and more. Markdown allows users to create structured documents using plain text, making it popular for documentation, readme files, and content creation. These Markdown elements allow for easy formatting of text and code, making documents more readable and structured while remaining simple to write and edit.

- Code Markdown: To format text as code in Markdown, you have two options:
 - Inline code: Use single backticks (`) to surround the code.
Example: `print("Hello, World!")`
 - Code blocks: Use triple backticks (```) before and after the code block, or indent each line with four spaces.
Example:
```  
def greet(name):  
 ...

```
 return f"Hello, {name}!"
```
```

Many Markdown processors also support syntax highlighting by specifying the language after the opening triple backticks:

```
```python  
def greet(name):
 return f"Hello, {name}!"
```
```

- Basic Text Formatting Markdown:
 - Bold: Use double asterisks or underscores.
Example: **bold text** or __bold text__
 - Italic: Use single asterisks or underscores.
Example: *italicized text* or _italicized text_
 - Strikethrough: Use double tildes.
Example: ~~strikethrough text~~
- Headers: Use hash symbols (#) at the start of the line. More hashes mean smaller headers.
Example:
Header 1
Header 2
Header 3
- Links: Use square brackets for the link text, followed by parentheses containing the URL.
Example: [OpenAI](https://www.openai.com)
- Lists: Use asterisks, plus signs, or hyphens for unordered lists; numbers for ordered lists.
Example:
 - * Item 1
 - * Item 2
 - * Subitem 2.1
 - 1. First item
 - 2. Second item
- Table Markdown

Table Markdown is a specific feature of Markdown that allows for the creation of tables using a simple text-based syntax. Here's a basic example:

```
```  
| Header 1 | Header 2 |
|-----|-----|
| Cell 1 | Cell 2 |
| Cell 3 | Cell 4 |
```
```

This syntax creates a table with headers and cells. The second line with dashes separates the header from the content. When rendered, this produces a neatly formatted table that's

both human-readable in its raw form and easily convertible to HTML or other formats.

Table Markdown is particularly useful for presenting structured data in a clear, organized manner without the need for complex formatting tools. It's widely supported in various Markdown processors and platforms, making it a versatile choice for data presentation in documentation, reports, and other text-based content.

Prompt Engineering Pipeline and Evaluation

Prompt engineering pipeline is a crucial aspect of developing effective AI systems. Let's explore several key aspects:

1. Test-Driven Prompt Engineering:

This approach applies the principles of test-driven development to prompt design. The process typically involves:

- Defining expected outcomes for specific inputs
- Writing prompts to achieve these outcomes
- Testing the prompts against a set of inputs
- Iteratively refining the prompts based on test results

This methodology helps ensure prompts are robust and perform consistently across various scenarios.

2. Evaluating Prompts:

Prompt evaluation is essential for measuring effectiveness and identifying areas for improvement.

Common evaluation metrics include:

- Accuracy: How often does the prompt produce the desired output?
- Consistency: Does the prompt generate similar results for similar inputs?
- Relevance: Are the outputs pertinent to the given task?
- Efficiency: How concise and targeted are the prompts in achieving the desired outcome?

3. Specialized Interfaces for Ground Truth Generation and Annotation:

Creating high-quality datasets for training and evaluation often requires specialized tools:

- Annotation platforms: Tools like LabelStudio or Prodigy can be adapted for prompt-specific annotation tasks.
- Crowdsourcing interfaces: Platforms like Amazon Mechanical Turk or Figure Eight can be used to gather diverse human-generated responses.
- Custom web applications: Tailored interfaces may be needed for complex or domain-specific annotation tasks.

4. Interfaces for Human Evaluation:

To assess prompt performance against ground truth, dedicated interfaces may be necessary:

- A/B testing platforms: Allow side-by-side comparison of different prompt outputs.
- Rating systems: Enable evaluators to score outputs on various dimensions (e.g., relevance, coherence, creativity).
- Qualitative feedback tools: Provide means for evaluators to give detailed comments and suggestions.

5. Bias Mitigation:

Addressing bias in prompt engineering is critical for creating fair and inclusive AI systems:

- Diverse input examples: Ensure prompts are tested with a wide range of inputs representing different demographics, perspectives, and scenarios.
- Bias detection tools: Utilize specialized software to identify potential biases in prompt outputs.
- Inclusive language guidelines: Develop and adhere to guidelines that promote neutral and inclusive language in prompts.
- Regular audits: Conduct periodic reviews to identify and address emerging biases.

6. Assimilating Human Feedback for Prompt Improvement:

Human feedback provides valuable insights that can significantly enhance prompt performance.

Here's how we can effectively incorporate this feedback:

- Structured Feedback Collection: Design specific rubrics for evaluators to provide targeted feedback on different aspects of prompt performance.
- Feedback Analysis: Use NLP techniques to extract key themes from qualitative comments and correlate feedback with prompt characteristics.
- Continuous Learning Loop: Implement a system for ongoing feedback collection and prompt refinement.
- Diverse Feedback Sources: Gather input from various stakeholders including end-users, domain experts, and diverse demographic groups.
- Balancing Act: Strike a balance between incorporating user preferences and maintaining the integrity of the AI system, while being mindful of potential biases in human feedback.

7. Drift Management:

As language models and their usage evolve, prompt effectiveness may change over time. Strategies to manage drift include:

- Continuous monitoring: Regularly assess prompt performance against established benchmarks.
- Version control: Maintain a history of prompt iterations to track changes and their impacts.
- Adaptive prompts: Design prompts that can accommodate shifting contexts or user behaviors.
- Feedback loops: Implement mechanisms to collect and incorporate user feedback for ongoing prompt refinement.

These aspects of prompt engineering are interconnected and often require an iterative approach. As the field evolves, new techniques and best practices continue to emerge, making it an exciting and dynamic area of AI development.

-----NOT part of the Exam

Pydantic Library

Popular frameworks for building LLM applications, such as LangChain and LlamaIndex, have integrated Pydantic as a core component of their parsing and data validation systems. This adoption underscores Pydantic's effectiveness and reliability in handling structured data from LLMs. For instance, LangChain extensively uses Pydantic in its output parsers, prompt templates, and chain inputs/outputs. The framework's PydanticOutputParser is a prime example, allowing developers to define expected output structures using Pydantic models. Similarly, LlamaIndex utilizes Pydantic for data validation and serialization in various components, including its document structures and query engines.

The widespread use of Pydantic in these frameworks not only validates its utility but also provides a consistent pattern for developers working across different LLM tools and libraries. With this context in mind, let's explore how Pydantic addresses the challenges of LLM output processing and how it can be integrated into your workflows.

What is Pydantic?

Pydantic is a Python library for data validation and settings management using Python type annotations. It enforces type hints at runtime and provides user-friendly errors when data is invalid. While not specifically designed for LLM output parsing, its features make it exceptionally well-suited for this task.

How Pydantic Addresses LLM Output Challenges

1. Structured Output Parsing:

Pydantic allows you to define data models as Python classes with type hints. These models serve as blueprints for your structured output. For example:

```
from pydantic import BaseModel

class Person(BaseModel):
    name: str
    age: int
    occupation: str

# Parsing LLM output
llm_output = '{"name": "John", "age": 30, "occupation": "Engineer"}'
person = Person.parse_raw(llm_output)
```

This approach provides a clear interface between the LLM output and your application logic.

2. Data Validation:

Pydantic automatically validates that the data conforms to the expected types and formats. It can handle complex validation rules, including field constraints, custom validators, and more.

3. Handling Complex Structures:

Unlike simpler parsing methods, Pydantic excels at handling nested data structures and lists of objects, making it suitable for more complex LLM outputs.

4. Error Handling:

When the LLM output doesn't match the expected format, Pydantic provides clear and informative error messages, making debugging easier.

5. Flexibility in Input Formats:

While Pydantic works well with JSON, it can also parse other formats. This flexibility aligns with the suggestion of using simplified output formats from LLMs.

Example: Parsing a Markdown Table with Pydantic

Combining the simplified output format suggestion with Pydantic, we might do something like this:

```
from pydantic import BaseModel
from typing import List

class Person(BaseModel):
    name: str
    age: int
    occupation: str

class PeopleTable(BaseModel):
    people: List[Person]

# LLM output as a markdown table
llm_output = """
Name	Age	Occupation
John	30	Engineer
Lisa	28	Designer
```

Parse the markdown table (simplified example)

```
lines = llm_output.strip().split('\n')[2:] # Skip header and separator
people = [Person(**dict(zip(['name', 'age', 'occupation'], line.split('|'))[:-1])) for line in lines]
```

```
1])) for line in lines]
table = PeopleTable(people=people)

print(table)
```

6. . Serialization:

Pydantic models can easily convert to and from JSON, dictionaries, and other formats, facilitating data exchange and storage.

7. Type Coercion and Data Fixing:

One of Pydantic's most powerful features is its ability to coerce and "fix" input data to match the defined types when possible. This makes it robust when dealing with real-world data that might not perfectly match your schema.

For example:

```
from pydantic import BaseModel
from typing import List, Optional

class User(BaseModel):
    id: int
    name: str
    email: str
    active: bool = True
    tags: Optional[List[str]] = None

class UserList(BaseModel):
    users: List[User]

user = User(id="1", name="Alice", email="alice@example.com", active="yes",
tags="python,coding")
print(user)
# Output: User(id=1, name='Alice', email='alice@example.com', active=True,
tags=['python', 'coding'])
```

In this case, Pydantic has:

- Coerced the string "1" to an integer for the `id` field
- Converted "yes" to `True` for the `active` field
- Split the comma-separated string into a list for the `tags` field

8. Strict Mode:

For cases where you want to disable automatic coercion, you can use `strict=True` on fields or entire models.

```
from pydantic import Field

class StrictUser(BaseModel):
    id: int = Field(..., strict=True)
```

By using Pydantic, you ensure that your output adheres to a specific structure while also benefiting from its data cleaning and coercion capabilities. This is particularly useful when:

- Working with APIs where you need to ensure consistent input and output structures
- Processing data from various sources that might have inconsistent formats
- Implementing data transfer objects (DTOs) in your application
- Creating configuration systems with validated settings

Pydantic's combination of strict typing with flexible data handling makes it an excellent tool for creating robust, well-defined structured outputs in Python applications.