

Revisiting Prompt Engineering: Strategic Interactions with Large Language Models for Application Development

Michael Borck¹

¹Business Information Systems, Curtin University, Perth Australia,

Corresponding author: Michael Borck, michael.borck@curtin.edu.au

Abstract

Prompt engineering has emerged as a crucial technique for leveraging the capabilities of large language models (LLMs) in application development. This paper explores the nuances of prompt engineering, presenting it not merely as a tool for interaction but as a programming paradigm that significantly eases the deployment of LLM capabilities into practical applications. We examine various levels and strategies of prompt engineering, highlighting their practical implications and potential for simplifying complex tasks.

Plain Language Summary

The paper discusses how prompt engineering can be used as a programming paradigm to effectively integrate large language models (LLMs) like GPT into software development. It describes two levels of prompt engineering: a basic level for straightforward tasks such as email drafting, and a more complex level involving programming interfaces for deeper software integration, such as automatic grading systems. The technique is portrayed as an empirical art that simplifies the development process by minimising coding efforts and maximising AI utility, although it faces challenges such as variability in effectiveness and potential high costs. The paper advocates for more research to enhance the scalability and efficiency of prompt engineering.

0.1 Introduction

The advent of large language models (LLMs) such as OpenAI’s GPT series has heralded a significant shift in the paradigms of programming and computational linguistics. Traditionally, the interaction with these models involved substantial coding efforts and a deep understanding of machine learning algorithms. However, a more accessible approach has surfaced, known as prompt engineering, which allows even those with minimal coding expertise to leverage the advanced capabilities of LLMs effectively. This paper delves into prompt engineering, defined as the strategic formulation of input prompts to maximise the utility of LLMs without internal model retraining. We argue that prompt engineering is an underappreciated yet powerful tool for developers, capable of achieving substantial outcomes with minimal effort.

Prompt engineering is defined as the strategic formulation of input queries to maximise the utility and accuracy of responses from LLMs. Initially perceived with skepticism, this methodology has quickly proven its worth, demonstrating that a well-crafted prompt can extract considerable value from these models with relatively low effort. The essence of prompt engineering lies not in rewriting the underlying algorithms of LLMs but in creatively interacting with them, using plain language to “program” or guide the models towards desired outputs.

The primary objective of this paper is to systematically explore and validate the effectiveness of various prompt engineering techniques. By comparing the performance of “good” versus “poor” prompts, this research aims to elucidate the impact of prompt construction on the quality of responses generated by LLMs. We employ a series of metrics, such as BERTscore and Perplexity, to quantitatively assess the responses, providing a rigorous analysis of how different prompting strategies affect the output of language models. This comparative study is intended to offer valuable insights for both developers and researchers, enabling them to harness the full potential of LLMs in diverse applications ranging from automated content generation to complex problem-solving tasks.

This paper includes embedded interactive elements that allow readers to directly engage with and replicate our computational analyses. By providing live data and code within the document, we aim to enhance transparency and applicability of our findings. It allows readers to not only digest the findings but also interact with the presented data and analyses, potentially replicating or building upon the research

in real-time. This integration promises to bridge the gap between theoretical exploration and practical application, making the research accessible and applicable to a broad audience interested in the future of artificial intelligence and machine learning.

0.2 Prompt Engineering: A Dual-Level Analysis

Prompt engineering operates at two distinct levels. The first, an intuitive and straightforward approach, utilises platforms like ChatGPT to execute simple tasks (e.g., text summarisation, email drafting). Although this method is accessible, its integration into larger systems is non-trivial. The second level involves a more sophisticated interaction through programming interfaces, such as Python APIs, enabling the integration of LLM responses into broader software ecosystems. This paper illustrates these levels with practical examples, including the development of an LLM-powered automatic grading system for educational applications.

0.3 Methodology: Empirical Techniques in Prompt Engineering

We categorise prompt engineering as an empirical art, where the composition and formatting of prompts are tailored to enhance model performance on specific tasks. This approach draws on the inherent desire of language models to complete text-based tasks, effectively “tricking” them into performing complex functions beyond simple text generation. We outline various techniques, such as the use of descriptive language, structured prompts, and the chaining of thought processes, to improve the interaction quality with LLMs.

This study systematically investigates the impact of prompt engineering on the performance of large language models (LLMs). Our approach involves comparing the efficacy of well-crafted (“good”) prompts against less optimised (“poor”) prompts in eliciting accurate and relevant responses from the models. We detail the experimental design, tools utilised, and metrics applied for evaluating the results.

0.3.0.1 Experimental Design

The experiments are structured to assess two main types of prompts:

1. **Good Prompts:** These are designed with specific linguistic and structural techniques believed to enhance the LLM’s understanding and response quality. Techniques include using descriptive language, structuring content, and employing contextual cues that guide the model towards the desired type of response.
2. **Poor Prompts:** These lack the strategic elements present in good prompts and typically include vague or ambiguous language that may lead to less accurate or relevant responses from the model.

For each type of prompt, identical tasks are presented to the LLM to ensure that the only variable affecting output quality is the prompt construction itself. Tasks include generating text-based responses across several domains, such as summarising passages, answering questions, and creating content based on given specifications.

0.3.0.2 Tools and Models

The study utilises OpenAI’s GPT-3 model due to its widespread adoption and robust performance across a variety of natural language processing tasks. We interact with GPT-3 using the following tools:

- **LangChain:** A Python library designed to facilitate the development of applications on top of LLMs. LangChain allows for structured interaction with GPT-3, managing prompt templates, and integrating response parsing mechanisms.
- **Quarto:** To document and share our findings, we use Quarto, an interactive computing framework that allows embedding live code, visualisations, and

narrative text. This choice supports the dynamic presentation of our methods and results, enabling readers to engage directly with the analyses.

0.3.0.3 Metrics for Evaluation

To objectively evaluate the quality of the responses generated by the LLM under different prompting conditions, we employ several metrics:

1. **BERTscore**: This metric computes the semantic similarity between the generated text and a set of reference texts. It is used to assess how well the content produced by the LLM aligns with expected outputs in terms of meaning and context.
2. **Perplexity**: Typically used to measure how well a probability model predicts a sample, perplexity in our context helps determine how “surprised” the model is by the responses it generates, which indirectly indicates the naturalness and fluency of the text.
3. **Manual Review**: To complement automated metrics, responses are also subjectively evaluated by human reviewers for relevance, coherence, and informativeness. This step ensures that our findings account for qualitative aspects of text generation that automated metrics may overlook.

0.3.0.4 Data Collection and Analysis

Responses from the LLM are collected under controlled conditions to ensure consistency across tests. Each response is logged along with the prompt used, allowing for detailed comparative analysis. The data are analysed using statistical tools to identify significant differences in performance metrics between good and poor prompts. Visualisations are created to illustrate these differences clearly, providing both quantitative and qualitative insights into the effectiveness of various prompt engineering techniques.

0.3.0.5 Comparative Approach

To evaluate the quality of prompts we treat the responses as references for each other. This can be an effective way to highlight which prompt elicits a more informative or relevant response on a given topic. Here’s how we implement this approach:

1. Define a topic or question you want to ask ChatGPT.
2. Craft two different prompts for the same topic - one that you consider a “poor” or suboptimal prompt, and one that you think is a “good” or well-designed prompt.
3. Provide both prompts to ChatGPT and obtain the responses.
4. Treat the response from the “good” prompt as the reference, and evaluate the response from the “poor” prompt against it using metrics like:
 - BERTScore: Calculate semantic similarity between the two responses using contextual embeddings.[1]
 - ROUGE/BLEU: Measure n-gram overlap between the responses.[2]
 - Human evaluation: Have human raters judge which response is more informative, relevant, and coherent.
5. Repeat the process, treating the response from the “poor” prompt as the reference, and evaluate the “good” prompt’s response against it.

By comparing the responses in this way, we identify which prompt leads to a more desirable output. A higher BERTScore, ROUGE, or BLEU score when using one response as the reference would indicate that the corresponding prompt produced a more relevant and informative response.[1][2][4]

Additionally, human evaluation can provide valuable qualitative insights into the strengths and weaknesses of each prompt, complementing the quantitative metrics.

This comparative approach leverages the fact that for the same topic, a well-crafted prompt should elicit a more relevant and high-quality response from the language model. By directly comparing the responses, you can empirically evaluate the effectiveness of your prompt engineering efforts.[1][2][4][5]

Citations: [1] <https://quaintitative.com/compare-llms/> [2] <https://andrewmaynard.net/comparative-prompts/> [3] <https://www.vcestudyguides.com/blog/the-five-types-of-text-response-prompts-archive> [4] <https://agio.com/comparing-ai-prompting-strategies/> [5] <https://typeset.io/questions/how-do-different-types-of-prompts-affect-the-quality-of-2rmbjrcisy>

Through this comprehensive methodology, the study aims to provide actionable insights into the art of prompt engineering, guiding users on how to best utilise LLMs for a range of applications.

Even without a predefined reference text, you can still compare the responses from a good prompt and a poor prompt using metrics like BERTScore or BLEU score. Here's how you can approach this:

1. Obtain responses from the language model using both the good prompt and the poor prompt.
2. Treat the response from the good prompt as the "reference" text and the response from the poor prompt as the "candidate" text.
3. Calculate the BERTScore or BLEU score between the reference (good prompt response) and candidate (poor prompt response).
4. Repeat the process by treating the poor prompt response as the reference and the good prompt response as the candidate.
5. Compare the scores in both directions to see which prompt elicited a response that is more similar/dissimilar to the other.

For BERTScore: - A higher BERTScore when using the good prompt response as reference indicates the poor prompt response is more semantically similar to the desired response.[1][4] - You can calculate BERTScore F1, precision, and recall to get a more nuanced comparison.

For BLEU score: - A higher BLEU score when using the good prompt response as reference suggests the poor prompt response has more n-gram overlap with the desired response.[5] - BLEU is based on precise word matching, so it may not fully capture semantic similarities.

This comparative approach lets you evaluate prompt quality without an external reference, by treating one prompt's response as the target output.[1][4][5]

Additionally, you can complement the metrics with human evaluation by having raters judge which response is more coherent, relevant and informative for the given topic.[1][4]

The key advantage is directly comparing the outputs from different prompts to assess which one yields a more desirable response from the language model.[1][4][5]

Citations: [1] <https://arxiv.org/pdf/2305.12421.pdf> [2] <https://www.vcestudyguides.com/blog/the-five-types-of-text-response-prompts-archive> [3] <https://quaintitative.com/compare-llms/> [4] <https://aclanthology.org/2021.wmt-1.59.pdf> [5] <https://aclanthology.org/P02-1040.pdf>

0.4 Results

This section presents the comprehensive results of our investigation into the effectiveness of prompt engineering techniques for enhancing the performance of large language models (LLMs). Our analysis focuses on quantitatively evaluating the impact of well-constructed versus poorly constructed prompts on the quality of the responses generated by LLMs. We utilise three established metrics for this evaluation: BLEU Score, ROUGE Score, and BERTscore, each providing insights into

different aspects of text quality such as precision, recall, and semantic similarity. The results clearly demonstrate the significant influence that prompt construction can have on the accuracy, relevance, and fluency of the generated text. Additionally, we provide detailed, interactive examples that not only illustrate these effects but also allow readers to explore the nuances of prompt engineering through live code. These examples showcase practical applications and underscore the practical implications of our findings, bridging theoretical research with actionable insights. Through dynamic visualisations and interactive Quarto cells, readers are invited to engage directly with the data, enhancing their understanding of how strategic prompt design can be effectively utilised in real-world applications.

0.4.1 Overview of Experimental Results

- Begin with a summary of the findings from the comparative analysis of good vs. poor prompts. This would involve presenting quantitative data from BERTscore, Perplexity, and manual reviews.
- Include charts, graphs, or tables that clearly depict the differences in LLM performance based on the type of prompt used.
- Compare the outputs and effectiveness of the LLM across different prompting techniques illustrated in the detailed examples and the case study. This section would analyse the broader implications of prompt engineering on practical applications.
- Discuss how variations in prompt construction can lead to significant differences in output quality and application functionality.

0.4.2 Detailed Examples of Good vs. Poor Prompts

This section provides detailed descriptions and Python code examples for various prompt engineering tricks used to enhance the performance of large language models (LLMs). Each strategy is designed to optimise the interaction with LLMs in different contexts, demonstrating practical applications and potential benefits.

0.4.2.1 Strategy 1: Be Descriptive (More is Better)

Explanation: Providing detailed information within prompts helps the LLM understand the context better, leading to more accurate and relevant responses. This strategy is particularly useful in scenarios where the details are critical to the task, such as generating personalised content or specific instructions.

Python Code Example:

```
# Importing the OpenAI GPT library
from openai import ChatCompletion

# Initialise the model with your API key
openai_api_key = 'your-api-key'
chat_model = ChatCompletion(api_key=openai_api_key, model="gpt-3.5-turbo")

# Define a detailed prompt for a birthday message
prompt = """
Write a birthday message for my dad. He is turning 60 years old, loves golf, enjoys classical music.
"""

# Generate the message using the detailed prompt
response = chat_model.create(
    messages=[{"role": "system", "content": prompt}]
)
print(response['choices'][0]['message']['content'])
```

0.4.2.2 Strategy 2: Give Examples (Few-Shot Learning)

Explanation: Demonstrating desired outputs through examples can significantly improve model performance, especially in tasks requiring specific formats or styles. This strategy is beneficial in educational settings, content creation, or anywhere model guidance through exemplars can enhance output relevancy.

Python Code Example:

```
# Define a prompt with examples for writing subtitles
prompt = """
Given the title of a blog post, write a subtitle that captures the essence of the article.

Title: Advances in Renewable Energy
Subtitle: Exploring the latest breakthroughs in sustainable power sources

Title: The Future of AI in Medicine
Subtitle: How AI is revolutionising diagnostics and patient care

Title: Innovations in Educational Technology
Subtitle: New tools that are transforming how students learn

Title: The Role of Genetics in Public Health
Subtitle:
"""

# Generate a subtitle for a new article
response = chat_model.create(
    messages=[{"role": "system", "content": prompt}]
)
print(response['choices'][0]['message']['content'])
```

0.4.2.3 Strategy 3: Use Structured Text

Explanation: Structuring prompts in a clear and organised manner can help the model parse and process the information more effectively. This approach is ideal for tasks that require data extraction, summarisation, or any application where clarity and precision are crucial.

Python Code Example:

```
# Define a structured prompt for a recipe
prompt = """
Create a recipe for a chocolate cake. The recipe should include:
- Title: Simple Chocolate Cake
- Ingredients: List all ingredients with precise measurements
- Instructions: Provide a step-by-step guide on how to prepare the cake
"""

# Generate the structured recipe
response = chat_model.create(
    messages=[{"role": "system", "content": prompt}]
)
print(response['choices'][0]['message']['content'])
```

0.4.2.4 Strategy 4: Chain of Thought

Explanation: Encouraging the model to “think” step-by-step can improve its reasoning abilities, particularly useful in problem-solving or tasks requiring logical progression, such as technical troubleshooting or complex decision-making.

Python Code Example:

```

# Define a prompt that uses the chain of thought approach
prompt = """
Problem: A user is unable to access their email account after multiple password resets. Let's t

Step 1: Verify if the user is using the correct email address.
Step 2: Check if their account is locked due to too many failed attempts.
Step 3: Ensure the password reset process is completed correctly.
Step 4: Suggest recovery through an alternate email or security questions.

What could be the issue based on these steps?
"""

# Generate a diagnostic response using the chain of thought
response = chat_model.create(
    messages=[{"role": "system", "content": prompt}]
)
print(response['choices'][0]['message']['content'])

```

256 *0.4.2.5 Strategy 5: Chatbot Personas*

257 **Explanation:** Prompting an LLM to adopt a specific persona can tailor its re-
 258 sponses to fit the desired character or expertise level. This strategy is particularly
 259 useful in customer service bots, educational tools, or any application where engaging
 260 and role-specific dialogue is beneficial.

261 **Python Code Example:**

```

# Importing the OpenAI GPT library
from openai import ChatCompletion

# Initialise the model with your API key
openai_api_key = 'your-api-key'
chat_model = ChatCompletion(api_key=openai_api_key, model="gpt-3.5-turbo")

# Define a prompt where the model adopts a persona
prompt = """
You are an expert gardener. A novice gardener asks for advice on starting a vegetable garden in
"""

# Generate a response using the gardener persona
response = chat_model.create(
    messages=[{"role": "system", "content": prompt}]
)
print(response['choices'][0]['message']['content'])

```

262 *0.4.2.6 Strategy 6: Flipped Approach*

263 **Explanation:** This technique involves prompting the LLM to ask questions back to
 264 the user to gain a better understanding of their needs. It's especially effective in con-
 265 sultation services or any scenario where clarifying user intent is crucial for accurate
 266 responses.

267 **Python Code Example:**

```

# Define a prompt that uses the flipped approach for better understanding
prompt = """
You are a travel consultant. A client is planning a trip to Europe but hasn't provided much det
"""

```



```
# Generate a series of questions to help tailor the travel advice
response = chat_model.create(
    messages=[{"role": "system", "content": prompt}]
)
print(response['choices'][0]['message']['content'])
```

0.4.2.7 Strategy 7: Reflect, Review, and Refine

Explanation: This strategy involves prompting the LLM to evaluate and potentially revise its previous responses. It is ideal for iterative tasks such as editing, programming, or any creative process where progressive refinement enhances the final output.

Python Code Example:

```
# Define a prompt that encourages reflection and refinement of a response
prompt = """
Initially, you wrote a brief summary of the novel 'To Kill a Mockingbird'. Now, review your summary and refine it to be more accurate and comprehensive.
"""

# Generate an improved summary through self-review and refinement
response = chat_model.create(
    messages=[{"role": "system", "content": prompt}]
)
print(response['choices'][0]['message']['content'])
```

Each of these tricks serves to enhance the interactivity and effectiveness of LLMs in specific scenarios. By employing these methods, developers can tailor the behavior of LLMs to meet diverse requirements and improve user experience, demonstrating the flexibility and potential of prompt engineering in practical applications. These examples illustrate how different prompt engineering tricks can be tailored to enhance LLM performance across a variety of tasks. The context in which each strategy is applied highlights its potential to improve the relevance and accuracy of the model's responses, making them invaluable tools in the arsenal of developers and researchers working with LLMs.

0.4.3 Comparative Analysis

0.4.3.1 BLUE Score

Here's an example code to compare the responses from a good prompt and a poor prompt using the BLEU score in Python:

```
from nltk.translate.bleu_score import sentence_bleu
import nltk

# Download required NLTK data
nltk.download('punkt')

# Define the good and poor prompts
good_prompt = "Explain the theory of relativity in simple terms."
poor_prompt = "What is relativity?"

# Get responses from the language model (replace with your own code)
good_response = "The theory of relativity, proposed by Albert Einstein, revolves around two main concepts: gravity and the speed of light."
poor_response = "Relativity is a scientific theory developed by Albert Einstein."

# Tokenise the responses
good_response_tokens = nltk.word_tokenize(good_response)
poor_response_tokens = nltk.word_tokenize(poor_response)
```

```

# Calculate BLEU scores
bleu_score_good_as_ref = sentence_bleu([good_response_tokens], poor_response_tokens)
bleu_score_poor_as_ref = sentence_bleu([poor_response_tokens], good_response_tokens)

print(f"BLEU score (good response as reference): {bleu_score_good_as_ref:.3f}")
print(f"BLEU score (poor response as reference): {bleu_score_poor_as_ref:.3f}")

```

Here's how the code works:

1. We import the necessary modules and download the required NLTK data for tokenisation.
2. We define the good and poor prompts as strings.
3. We obtain the responses from the language model for the good and poor prompts. In this example, we've hardcoded the responses, but in practice, you would replace these with your own code to generate responses from the language model.
4. We tokenise the good and poor responses using `nltk.word_tokenize`.
5. We calculate the BLEU scores in both directions:
 - `bleu_score_good_as_ref` treats the good prompt response as the reference and the poor prompt response as the candidate.
 - `bleu_score_poor_as_ref` treats the poor prompt response as the reference and the good prompt response as the candidate.
6. We print out both BLEU scores.

When you run this code, you should see output similar to:

```

BLEU score (good response as reference): 0.235
BLEU score (poor response as reference): 0.087

```

In this example, the BLEU score is higher when using the good prompt response as the reference (0.235) compared to using the poor prompt response as the reference (0.087). This suggests that the response from the good prompt has more n-gram overlap with the response from the poor prompt, indicating that the good prompt elicited a more informative and relevant response.

You can interpret the BLEU scores as follows:

- A higher BLEU score when using the good prompt response as the reference suggests that the poor prompt response is more similar to the desired response (from the good prompt).
- A lower BLEU score when using the poor prompt response as the reference indicates that the good prompt response is less similar to the undesirable response (from the poor prompt).

By comparing the BLEU scores in both directions, you can evaluate which prompt yielded a response that is more relevant and informative for the given topic.

0.4.3.2 ROUGE (Recall-Oriented Understudy for Gisting Evaluation)

The ROUGE score is a set of metrics used primarily for evaluating automatic summarisation and machine translation. It measures the overlap between the system-generated text and reference texts, focusing on the number of overlapping units such as n-grams, word sequences, and word pairs. The most commonly used variants are ROUGE-N (which measures n-gram overlap), ROUGE-L (which measures the longest common subsequence), and ROUGE-S (which measures skip-bigram overlap). Below, we implement the ROUGE-L score in Python to evaluate the quality of responses from language models. This will require the installation of the `rouge` package, which can be done using `pip`:

```
pip install rouge
```

0.4.4 Python Implementation Example for ROUGE-L

Here's an example of how you can calculate the ROUGE-L score for a language model's responses using the Python `rouge` library:

```
from rouge import Rouge

# Initialise the ROUGE scorer
rouge = Rouge()

# Example reference and system-generated summaries
reference = "John F. Kennedy was the 35th president of the United States. He served from 1961 u
system_generated_1 = "JFK, also known as John Kennedy, was the president of the US who served f
system_generated_2 = "The 35th president of the United States was John F. Kennedy, who served f

# Calculate ROUGE scores
scores_1 = rouge.get_scores(system_generated_1, reference)
scores_2 = rouge.get_scores(system_generated_2, reference)

# Print ROUGE-L scores for each system-generated summary
print("ROUGE-L Score for System-Generated Summary 1:", scores_1[0]['rouge-l'])
print("ROUGE-L Score for System-Generated Summary 2:", scores_2[0]['rouge-l'])
```

0.4.5 Explanation

- **ROUGE-L:** This score focuses on the longest common subsequence between the system-generated text and the reference text. It considers sentence-level structure similarity naturally and identifies longest co-occurring in-sequence n-grams of words. It provides two scores: precision (what fraction of the generated words are relevant) and recall (what fraction of the reference words were captured by the generated text), and an F-measure which is the harmonic mean of precision and recall.
- **System-Generated Summaries:** These are the texts generated by your system, in this case, responses from a language model.
- **Reference Summary:** This is the “ideal” or target summary against which the system-generated summaries are compared.

0.4.6 Usage Considerations

- **Contextual Relevance:** While ROUGE scores can provide objective measures of textual overlap, they might not fully account for the semantic accuracy or contextual relevance of the generated text. For more nuanced evaluations, consider combining ROUGE with manual qualitative assessments.
- **Variants:** Depending on the specific requirements of your evaluation, you might choose different ROUGE metrics (like ROUGE-N for n-gram overlap, or ROUGE-S for skip-bigrams) to focus on different aspects of the comparison.

This approach allows you to quantitatively assess how well a language model's responses match a reference, providing a standardised measure to compare different prompting strategies or model configurations in your research or application development.

0.4.7 Using BERTscore without a Predefined Reference

Using BERTscore to compare the responses to a “poor” prompt and a “good” prompt is a viable approach, especially when you are trying to assess how effectively each prompt elicits the intended response from a language model. Essentially, you would compare the response to each type of prompt against a reference standard

or even against each other to determine which prompt results in a more semantically rich and relevant response.

If you don't have a predefined reference response but want to compare the quality of responses between two prompts, you can consider the following approaches:

1. **Use Responses as Mutual References:** Treat the response from one prompt as the reference for the other and vice versa. This comparative approach can highlight which prompt elicits a response that is more informative or relevant to the topic.
2. **Create a Synthetic Reference:** If you know what information the response should contain, you can construct a synthetic reference response that includes all the expected elements. This reference can then be used to evaluate the responses from both the poor and good prompts.

Here's how to implement the first approach using the BERTscore:

0.4.7.1 Implementation of BERTscore Comparison

First, ensure you have the `bert-score` package installed:

```
pip install bert-score
```

Then, set up your comparison:

```
from bert_score import score

# Assuming you've obtained responses from the model
response_from_poor_prompt = "The quick brown fox tries to jump but fails."
response_from_good_prompt = "The quick brown fox jumps over the lasy dog."

# Using each other as references
P_poor, R_poor, F1_poor = score([response_from_poor_prompt], [response_from_good_prompt], lang=
P_good, R_good, F1_good = score([response_from_good_prompt], [response_from_poor_prompt], lang=

# Print results
print("Comparison using Good Prompt as Reference:")
print(f"Precision: {P_poor.tolist()}, Recall: {R_poor.tolist()}, F1 Score: {F1_poor.tolist()}")

print("\nComparison using Poor Prompt as Reference:")
print(f"Precision: {P_good.tolist()}, Recall: {R_good.tolist()}, F1 Score: {F1_good.tolist()}")
```

0.4.8 Interpretation

- **Precision:** Measures how much of the information in the generated response is relevant (i.e., how much of the generated content actually pertains to what the reference response discusses).
- **Recall:** Measures how much of the reference's content is covered by the generated response (i.e., how much of the essential information in the reference was captured in the generated response).
- **F1 Score:** The harmonic mean of precision and recall, providing a single score that balances both completeness and accuracy.

0.4.9 Considerations

This method of using responses as mutual references can be particularly useful when direct comparison metrics are needed to evaluate the effectiveness of different prompts in eliciting detailed and relevant responses. It is somewhat subjective, as it assumes that one of the responses contains sufficient quality content to serve as a benchmark—a fair assumption if one prompt is indeed superior.

By comparing these metrics, you can objectively analyse which prompt results in a better response in terms of information richness and relevance, thereby supporting the effectiveness of your prompt engineering techniques.

0.5 Case Study: Building an Automatic Grader

The practical application of prompt engineering is demonstrated through the development of an automatic grader for educational purposes. Using LangChain, a Python library for LLM integration, we construct a prompt-based system that evaluates student responses to open-ended questions. This system exemplifies the shift from traditional programming methods to a model where significant portions of logic are outsourced to an LLM, highlighting the efficiency and scalability of prompt engineering.

This section provides a detailed guide to building an automatic grader application utilising LangChain, a Python library designed to simplify the integration of large language models (LLMs) in application development. The application leverages prompt engineering to evaluate student responses in a high school history class, demonstrating the practical implementation of LLMs in educational settings.

0.5.0.1 Overview

The automatic grader application is designed to assess short-answer questions, where answers might vary but still be correct. The grader needs to handle different phrasings, synonyms, and minor spelling errors effectively.

0.5.0.2 Setup and Dependencies

Before starting, ensure you have Python installed along with the necessary libraries. You'll need `langchain` and `openai`. If not already installed, they can be added via `pip`:

```
pip install langchain openai
```

0.5.0.3 Code Implementation

Step 1: Import Necessary Libraries

```
# Import necessary modules from LangChain and OpenAI
from langchain.chat_models import ChatOpenAI
from langchain.prompts import PromptTemplate
from langchain.chains import LLMChain
from langchain.schema import BaseOutputParser
```

Step 2: Define the Prompt Template The prompt template mimics a high school history teacher grading an assignment. It should account for the question, the correct answer, and the student's response.

```
# Define the prompt template for grading
prompt_template_text = """
You are a high school history teacher grading homework assignments.
Based on the homework question indicated by "**Q:**" and the correct answer
indicated by "**A:**", your task is to determine whether the student's answer is correct.
Grading is binary; therefore, student answers can be correct or wrong.
Simple misspellings are okay.

**Q:** {question}
**A:** {correct_answer}

**Student's Answer:** {student_answer}
"""
```

Step 3: Initialise LangChain

```
# Define and initialise the ChatOpenAI model with your OpenAI API key
chat_model = ChatOpenAI(openai_api_key="your-openai-api-key", temperature=0)

# Create a prompt template object
prompt = PromptTemplate(
    input_variables=["question", "correct_answer", "student_answer"],
    template=prompt_template_text
)

# Define the chain using the LLM and the prompt
chain = LLMChain(llm=chat_model, prompt=prompt)
```

Step 4: Output Parser This component converts the LLM's response into a structured format that can be easily interpreted by other systems.

```
# Define an output parser to interpret the LLM's grading response
class GradeOutputParser(BaseOutputParser):
    """Parses the LLM's response to determine if the answer is correct or wrong."""

    def parse(self, text: str):
        """Check if the response indicates the student's answer was wrong."""
        return "wrong" not in text.lower()

# Update the chain to use the new output parser
chain.output_parser = GradeOutputParser()
```

Step 5: Running the Grader

```
# Sample question and answers
question = "Who was the 35th president of the United States of America?"
correct_answer = "John F. Kennedy"
student_answers = ["JFK", "John F Kennedy", "John Kennedy", "Jack Kennedy"]

# Evaluate each student answer
for student_answer in student_answers:
    result = chain.run({'question': question, 'correct_answer': correct_answer, 'student_answer': student_answer})
    print(f"Question: {question}")
    print(f"Student Answer: {student_answer} - {'Correct' if result else 'Incorrect'}\n")
```

This example application showcases how LangChain can be utilised to develop practical, LLM-integrated solutions for real-world problems such as grading. The automatic grader is not only efficient but also demonstrates the capability of LLMs to handle variability in natural language processing tasks, making it an ideal solution for educational applications.

0.6 Discussion: Efficacy and Limitations of Prompt Engineering

While prompt engineering offers substantial advantages, including significant reductions in development time and increased flexibility, it is not without limitations. The effectiveness of prompt strategies can vary significantly between different LLM versions, and the approach may incur considerable computational costs. Furthermore, the general-purpose nature of large models like ChatGPT may not be optimal for specialised tasks, which could be better served by fine-tuning specific models tailored to particular needs.

0.7 Conclusion

Prompt engineering represents a transformative approach to programming, enabling developers to harness the capabilities of LLMs in a user-friendly and cost-effective

manner. As this field evolves, it is expected that the techniques and strategies of prompt engineering will become increasingly refined, paving the way for more sophisticated and integrated applications across various domains.

0.8 Future Work

Further research is needed to explore the scalability of prompt engineering across different platforms and its integration with advanced model fine-tuning techniques. Additionally, comparative studies could elucidate the conditions under which prompt engineering is most effective compared to traditional programming approaches.

0.8.1 References

1. Talebi, S. (2023). Prompt Engineering: How to strategy AI into Solving Your Problems. *Towards Data Science*.
2. Karpathy, A. (n.d.). *Language Models as Few-Shot Learners*.
3. Wei, J., et al. (2021). *Chain of Thought Prompting Elicits Reasoning in Large Language Models*.