

11-411/11-611 Homework Assignment 4: Entity Recognition

Lara Marinov - Spring 2024

Abstract

In this report, I experiment with multiple different methods of constructing demonstrations and prompts for the OpenAI large language model (LLM) to try to improve its performance on token classification. I first experimented with changing the number of shots on a baseline model, trying to discover which value of shots performed best. I then experimented with using three LLMs - Perplexity AI, ChatGPT 3.5, and Gemini - to create modified prompts for the OpenAI LLM. I conclude that across both baseline approaches and approaches with modified prompts, using 40 shots and the baseline prompt worked best with an F1 score of 0.361. Of the three modified prompt approaches, ChatGPT's "fill-in-the-blank" prompt style worked best with an F1 score of 0.302.

Part 1: Methodology

Background and Justification

Prompt engineering has been an active area of NLP research, focusing on techniques to improve the performance of large language models (LLMs) on specific tasks by carefully crafting prompts. Many formatting and high-level techniques for prompt engineering have been previously explored like few-shot prompting, where the prompt includes examples of the desired task to guide the LLM generation [The New Stack \(2023\)](#), and chain-of-thought prompting, where the LLM is instructed to break down its reasoning process into multiple steps [Chen et al. \(2023\)](#).

Along with high level formatting techniques, other research has been done on smaller scale changes which also improve LLM performance such as using affirmative directives, incorporating examples, assigning roles to LLMs, and using output primers [Bsharat et al. \(2024\)](#).

The techniques mentioned above have shown promise in improving LLM performance, but they primarily focus on generating highly likely or factual outputs. There has also been research into generating relevant, but less probable prompts, which can be valuable in decision-making scenarios where considering alternative possibilities is important [Tang et al. \(2023\)](#).

While significant research has explored multiple prompt engineering techniques, little research exists which has explored utilizing LLMs themselves as auxiliary tools within the overall prompt engineering process. Traditionally, research has focused on improving and refining human-crafted prompts for direct LLM interaction. This short report aims to investigate the potential of LLMs to contribute to the development of effective prompts for other LLMs.

Section 1.1: Baseline Experimentation

To start, I prompted the OpenAI LLM with the baseline format given to us in the notebook which can be seen in [Figure 1](#). By varying only the number of shots and keeping the same prompt structure, I could understand the effect of training data size on the model's performance in an isolated environment. This allowed me to assess how OpenAI LLM's ability to identify and classify entities changed as the amount of labeled examples provided increased. My subsequent runs then introduced variations to the prompt structure itself, building upon the basic understanding of the model's behavior using the baseline format. The initial evaluation served as a sort of control experiment, establishing a benchmark for how the OpenAI LLM performed with a fixed prompt format before introducing additional complexities.

The LLM performed reasonably well, with the F1 score increasing as the number of shots increased. The highest F1 score was 0.361, which was achieved using a 40-shot approach. More results can be seen in [Figure 4](#) in Section 2. I wanted to run additional experiments, increasing the number of shots to 100, but I also recognized that it was important to experiment with different prompt structures before I ran out of credits. Given the resource limitation and lack of information regarding how many credits I had left, I decided to switch focus to changing the prompt structure and revisit the impact of higher shot sizes if I still had credits. I unfortunately ran out of credits before I could return to run with increased shot values. [Section 3.4](#) discusses the limitations encountered and how they affected the rest of my report in more detail. With more credits, I would have increased the shot size to at least 100, and if performance continued to improve, I may have decided to further increase the shot size to see how performance would change.

```

1 messages = [
2     {'role': SYSTEM_STR, MSG_STR:
3         ""You will be given input text containing different types of entities that you will
4         label.
5         This is the list of entity types to label: Deity, Mythological_king,
6         Cretaceous_dinosaur, Aquatic_mammal, Aquatic_animal, Goddess.
7         Label the entities by surrounding them with tags like '<Cretaceous_dinosaur>
8         Beipiaognathus </Cretaceous_dinosaur>'.""
9     },
10    {'role': USER_STR, MSG_STR: ""Text: Once paired in later myths with her Titan brother
    Hyperion as her husband, mild-eyed Euryphaessa, the far-shining one of the Homeric Hymn
    to Helios, was said to be the mother of Helios (the Sun), Selene (the Moon), and Eos (
    the Dawn).""},
11    {'role': SYSTEM_STR, MSG_STR: ""Labels: Once paired in later myths with her Titan
    brother <Deity> Hyperion </Deity> as her husband, mild-eyed Euryphaessa, the far-shining
    one of the Homeric Hymn to Helios, was said to be the mother of Helios (the Sun), <
    Goddess> Selene </Goddess> (the Moon), and <Goddess> Eos </Goddess> (the Dawn).""},
12    {'role': USER_STR, MSG_STR: ""Text: From her ideological conception, Taweret was
    closely grouped with (and is often indistinguishable from) several other protective
    hippopotamus goddesses: Ipet, Reret, and Hedjet.\nLabels: ""}]

```

Figure 1: The baseline prompt format provided in the notebook

Section 1.2: Modified Prompts Experimentation

In an effort to explore a unique prompt engineering technique, I investigated the feasibility of utilizing an LLM to generate prompt structures for the OpenAI LLM backend. This approach involved getting prompt structures from three well-known conversational LLMs: Perplexity AI, ChatGPT, and Gemini. All three models re-

ceived the exact same sequence of two prompts. The first prompt, seen in **Figure 2**, outlined the task and the desired entity tags for classification. I included the baseline prompt into the message so that the models could understand the context of the current prompts.

```

I am working on prompt engineering for an LLM which is doing token
classification. This is the current structure of the prompt:

chat.history = [{'role': SYSTEM_STR, MSG_STR: ""You will be given
input text containing different types of entities that you will label.
This is the list of entity types to label: Deity, Mythological_king,
Cretaceous_dinosaur, Aquatic_mammal, Aquatic_animal, Goddess. Label
the entities by surrounding them with tags like '<Cretaceous_dinosaur>
Beipiaognathus </Cretaceous_dinosaur>'.""},
{'role': USER_STR, MSG_STR: ""Text: Once paired in later myths with her
Titan brother Hyperion as her husband, mild-eyed Euryphaessa, the far-shining
one of the Homeric Hymn to Helios, was said to be the mother of Helios (the
Sun), Selene (the Moon), and Eos (the Dawn).""},
{'role': SYSTEM_STR, MSG_STR: ""Labels: Once paired in later myths with
her Titan brother <Deity> Hyperion </Deity> as her husband, mild-eyed
Euryphaessa, the far-shining one of the Homeric Hymn to Helios, was said
to be the mother of Helios (the Sun), <Goddess> Selene </Goddess> (the Moon),
and <Goddess> Eos </Goddess> (the Dawn).""}]

Write a few options for different prompts I could give to the LLM

```

Figure 2: The prompt given to the three LLMs

Each LLM gave 3-4 suggestions for approaches with example messages. All three responses shared some similarities in the approaches they suggested, but there were also considerable differences. Each LLM suggested basic methods such as zero shot and few shot prompting, which I already covered in the experiments with the baseline prompts. To decide which prompt structure suggestion I would use from each model, I chose the suggestion most dissimilar to the responses from the other LLMs. I then prompted each LLM once more with the message in **Figure 3**, only switching out the text with {“examples with explanations”, “fill in the blank”, ‘error-augmented prompts’} for the respective LLM.

Section 1.2.1: Perplexity AI – Examples with explanations

The “examples with explanations” prompt engineering technique is a powerful, intuitive approach where examples are provided along with correct answers and explanations. The core idea is to influence the LLM by showing it examples of the desired task along with the correct outputs and, crucially, explanations for why those outputs are correct [Weng \(2023\)](#). This technique helps the LLM better understand the task and the reasoning required to generate the desired classification. By seeing the explanations of the correct answers, the model can learn the patterns, logic, and context behind the inputs and can apply that knowledge to future inputs.

For the specific entity recognition task of this assignment, the input text and labels would be the same as in the baseline approach. The explanation added to each pair of input text and labels would include a description of why the specific special tags were chosen for each token. For example, a possible explanation could describe that “Athena” is a Greek goddess which is why she is tagged with the “Goddess” tag. **Appendix Section A1** shows the complete output from Perplexity AI which was used as the precedent to create the rest of the message triples sent to the OpenAI backend.

Section 1.2.2: ChatGPT – Fill in the blank

The “fill-in-the-blank” prompt engineering technique for token classification is a way to structure input prompts so that the LLM is asked to identify and label specific entities within a given text by replacing blank spaces with

the appropriate tags [Liu et al. \(2021\)](#).

Instead of providing the labeled text associated with each input directly, the prompt presents the text with blank spaces where the HTML tags should be inserted: `< __ > Hyperion < __ >`.

By presenting the task in this “fill-in-the-blank” format, the LLM is prompted to identify the relevant entities within the text and replace the blank spaces with the appropriate tags. One significant potential advantage of this structure is that the LLM, instead of passively receiving the labeled text, is actively engaged in the task of identifying and labeling the tokens. The hope is that this active involvement leads the LLM to better understand and remember the context and goal of the task.

There is also an extension of the “fill-in-the-blank” format where, in addition to examples with blank tags, a few examples of correctly labeled texts are also included. These correct examples could possibly improve the LLM’s performance by helping it learn the desired output format of the entity labeling.

Section 1.2.3: Gemini – Error-augmented prompts

The “error-augmented prompts” prompt engineering technique for token classification involves intentionally introducing errors, such as filler words, spelling mistakes, or inconsistent capitalization, into the prompts provided to the LLM.

By including different types of errors in the input messages, the model is forced to learn to identify and extract the relevant entities and their tags despite distractions and irregularities being present in the input. This can be particularly useful for real-world applications where the input text may not always be perfectly formatted or error-free [Eliot \(2023\)](#).

By training the model on prompts with various types of errors, it can learn to extract relevant information better and it can handle a wider range of inputs. The hope is that because the model becomes less sensitive to minor formatting issues, typos, or inconsistencies in the input text, it is able to focus its attention on the accuracy of the entity tags. The model may develop a better understanding of the context, which would further help it differentiate relevant entities.

Expand more on the <see above> approach and give an example prompt in code using a structure like the one I provided in my previous message.

Figure 3: The follow-up prompt given to each LLM to receive uniform prompts to feed to the OpenAI backend

Part 2: Experimental Results

Approach	Shots	Accuracy	Precision	Recall	F1
Finetune BERT	-	0.948	0.410	0.509	0.454
Zero-shot LLM Baseline	0	0.965	0.460	0.239	0.296
One-shot LLM Baseline	1	0.959	0.322	0.300	0.311
Few-shot LLM Baseline	5	0.965	0.389	0.239	0.296
Few-shot LLM Baseline	10	0.965	0.329	0.212	0.258
Few-shot LLM Baseline	20	0.948	0.366	0.258	0.302
Few-shot LLM Baseline	40	0.946	0.441	0.306	0.361
Perplexity AI suggestion	5	0.964	0.363	0.242	0.298
ChatGPT suggestion	5	0.956	0.370	0.256	0.302
Gemini suggestion	5	0.958	0.347	0.263	0.290

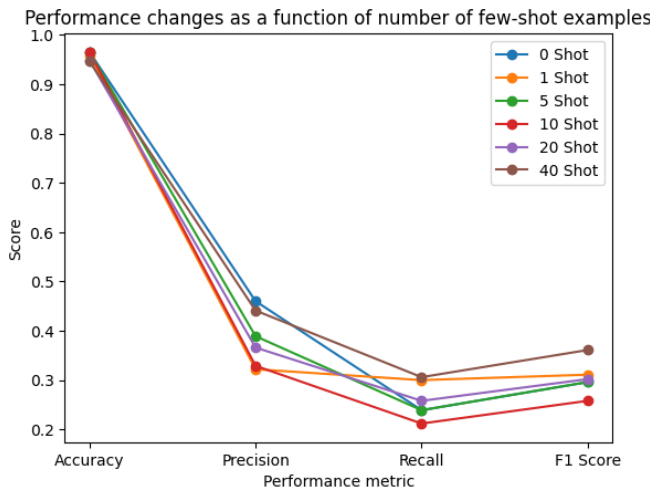
Figure 4: Experimental results table, best result (excluding BERT) in each column is shown in **bold**

As discussed in **Section 1**, I chose to experiment with the feasibility of utilizing an LLM to generate prompt structures for the OpenAI LLM for the token classification task. My experimental results table is displayed above. The first observation to note is that the approach with the highest F1 score was using the LLM baseline with 40 shots. This result is not surprising considering the nature of the specific labeling we are asking the model to perform. Although, some discussion suggests that larger models do not need as many shots because they have seen more data in training and can thus generalize better

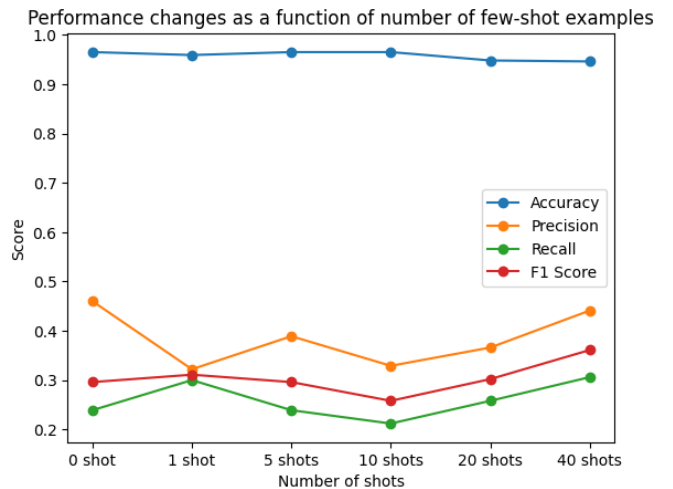
to a wide variety of tasks (Reynolds & McDonell (2021); Weird_Foundation5044 (2023)).

The second interesting result to note is that the accuracy across all runs is far higher than the score on any other metric, regardless of approach. This result is discussed in depth in **Section 3** of this report. All runs with modified prompt structures were done using 5 shots. More runs were planned with increased shots, but the restricted number of API credits limited the number of runs I was able to accomplish. These limitations are discussed further in **Section 3.4**.

Performance changes



(a) Performance changes plotted by metric



(b) Performances changes plotted by shots

Figure 5: Performances changes with baseline prompt structure and changes to number of shots

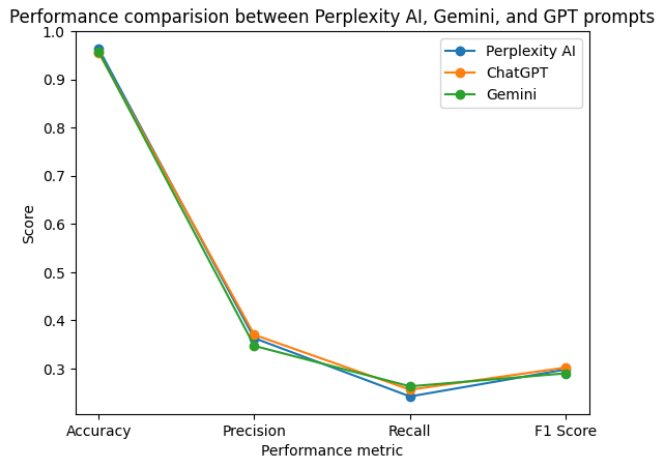


Figure 6: Performances changes plotted by metric

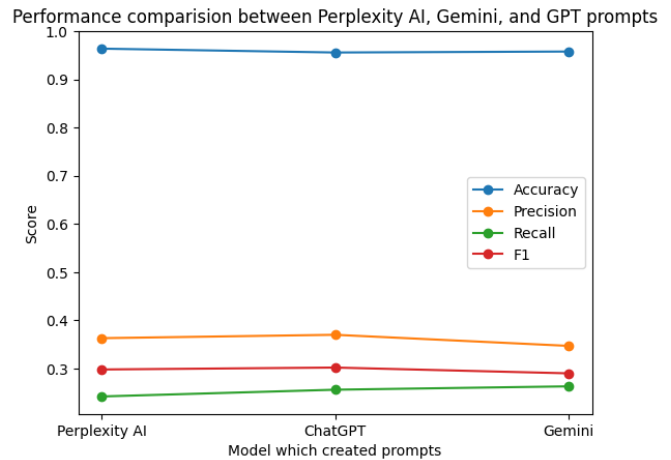


Figure 7: Performances changes plotted by run type

Figure 8: Performances changes with prompt structure modification

Figure 5 shows how the performance changes as a function of number of few-shot examples. Once again, accuracy is significantly higher across all of the shots, as can be seen in **Figure 5b**. The run with 40 shots performed better across all four metrics, although the zero-shot run also had very similar precision to the 40-shot run. **Figures 6 - 8** show the same performance changes, but on the runs with modified prompt structures.

The performances across the different prompt structures in comparison to each other are more similar than the performances with different shots in comparison to each other. The “fill in the blank” prompt style from ChatGPT marginally outperformed the two styles from Perplexity AI and Gemini. These results are discussed in more depth in **Section 3**.

Part 3: Analysis and Discussion

		'O'	'Aquatic_animal'	'Aquatic_mammal'	'Cretaceous_dinosaur'	'Deity'	'Goddess'	'Mythological_king'
Train set	Raw Number	161563	1304	714	450	1092	997	313
	Raw Number (B I)	–	739 565	451 263	438 12	969 123	852 145	247 66
	Fraction of total	0.9699	0.0081	0.0044	0.0028	0.0068	0.0062	0.0019
	Fraction of category (B I)	–	0.5667 0.4333	0.6317 0.3683	0.9733 0.0267	0.8874 0.1126	0.8546 0.1454	0.7891 0.2109
Dev set	Raw Number	13311	101	49	39	128	78	16
	Raw Number (B I)	–	60 41	34 15	36 3	114 14	68 10	14 2
	Fraction of total	0.9700	0.0074	0.0036	0.0028	0.0093	0.0057	0.0012
	Fraction of category (B I)	–	0.5941 0.4059	0.6939 0.3061	0.9231 0.0833	0.8906 0.1094	0.8718 0.1282	0.875 0.125
Both sets	Raw Number	170004	1405	763	489	1220	1075	329
	Raw Number (B I)	–	799 606	485 278	474 15	1083 137	920 155	261 68
	Fraction of total	0.9699	0.0080	0.0044	0.0028	0.0069	0.0061	0.0019
	Fraction of category (B I)	–	0.5687 0.4313	0.6356 0.3644	0.9693 0.0307	0.8877 0.1123	0.8558 0.1442	0.7933 0.2067

Figure 9: Breakdown of entity type tag distribution across dev and training datasets

Section 3.1 Analysis The analysis of the experimental results revealed that while the performance improvement for modifications to the prompt structure were relatively minor, increasing the number of shots consistently improved model performance. This suggests that the model benefits from a larger dataset, even with marginally different prompts.

However, analyzing the breakdown of tags in both datasets in **Figure 9** and **Figures 10 - 12** revealed that the datasets are overwhelmingly skewed towards the 'O' tag. As observed in **Section 2**, the accuracy across all runs – both more shots and modified prompts – was significantly higher than the rest of the metrics, which is likely a direct cause of this imbalance; even the most simplistic model that predicts only 'O' for every token could achieve a deceptively high accuracy score.

This observation also raised the question: how do the modified prompt structures possibly reduce this tag imbalance problem? One plausible explanation is that the different prompt structures can help the model learn that the task is not just about identifying each token as the 'O' tag but also recognizing and labeling other entity types based on contextual information. By exposing the model to various prompt structures with examples and explanations, it may capture the nuances and patterns required for accurate tag recognition better, rather than defaulting to a naive 'O' prediction strategy.

To further investigate the effectiveness of the modified prompts, I thought it would be insightful to analyze the performance metrics across the different prompt formats, broken down by entity tag type. I thought that this granular analysis could reveal valuable insights into

which prompting techniques are more adept at identifying specific entity types, rather than simply relying on the majority 'O' class predictions consistently.

Figures 10 - 12 shows the breakdown of the tags in the dev and training sets, excluding 'O'. Clearly, some labels are more common than others: the 'Deity' and 'Goddess' tags appear more frequently than tags like 'Mythological_king'. So, if a model is good at predicting the 'Deity' and 'Goddess' tags, we would expect it to perform better overall, even if it's worse at predicting many of the other tags. This suggestion is supported by the performance of the 40-shot run, which had the highest F1 score on the 'Deity' and 'Goddess' tags (as seen in **Figures 16 and 17**) and the highest overall F1 score, even though it was not consistently the best at predicting every tag. These results could suggest that to enhance overall model performance in the future, prompts could strategically emphasize the less frequently encountered tags. This approach would aim to address the potential bias towards frequent tags and it would encourage the model to focus on tags where its learning can be further optimized.

Figures 13 - 18 also reveal that there was not one run type which had better performance across all tags. Indeed, quite the opposite is true as each of the four run types: 40-shot, Perplexity AI, ChatGPT, and Gemini had at least one tag for which they performed the best (as measured by the F1 score). This suggests that combining the prompting style used in the Perplexity AI, ChatGPT, and Gemini runs and increasing the numbers of shots (possibly to even more than 40) could further improve the model's performance.

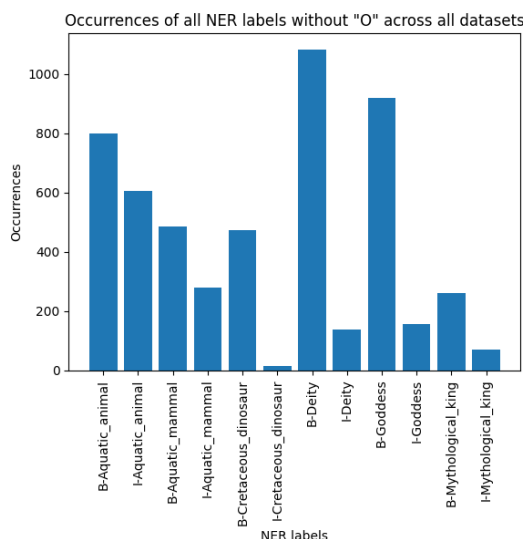


Figure 10: Distribution of tags across both the dev and training sets; generation code in Appendix Section A5

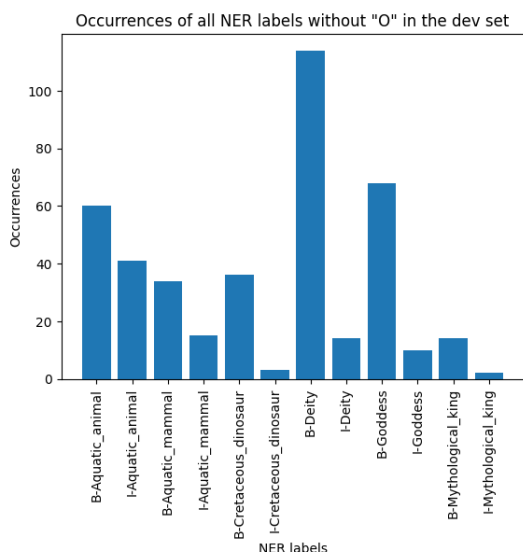


Figure 11: Distribution across the dev set; generation code in Appendix Section A5

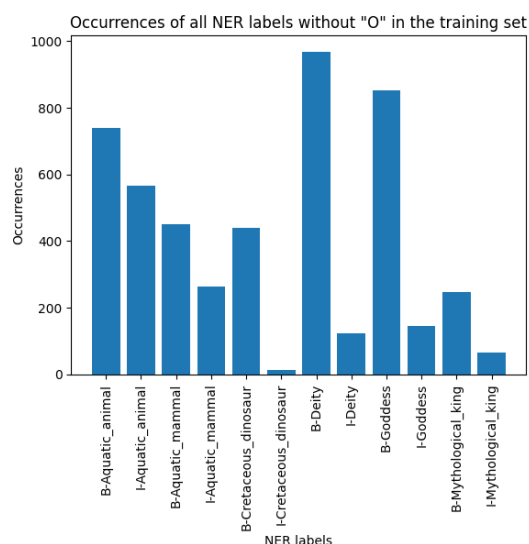


Figure 12: Distribution across the training set; generation code in Appendix Section A5

Further, if a particular prompt structure consistently outperforms others in recognizing and labeling rare entity types, such as “Mythological.king” or “Cretaceous.dinosaur,” it could indicate that the prompt structure is better at capturing the contextual cues and patterns associated with those entities. Identifying the prompt structures which result in better performance on rare tags is crucial. These prompt structures could be interwoven with ones that perform better on common tags, possibly resulting in models that predict both common and rare tags more accurately.

Section 3.2 Future plans It would be interesting to

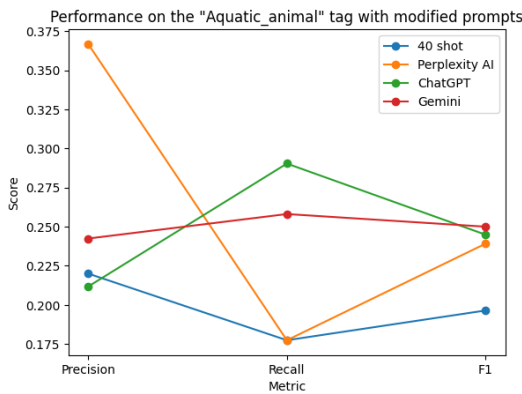
further investigate the potential of leveraging LLMs as auxiliary prompt engineering tools. The research could explore the efficacy of LLMs in generating effective prompt structures along with compiling a comparative analysis of the strengths and weaknesses of these LLM-generated prompts.

Section 3.3 Conclusion In this short report, I explored different methods of constructing demonstrations and prompts for the OpenAI LLM to try to improve its performance on token classification. I experimented both with a baseline prompt structure, varying only the number of shots, and with three different modified prompt

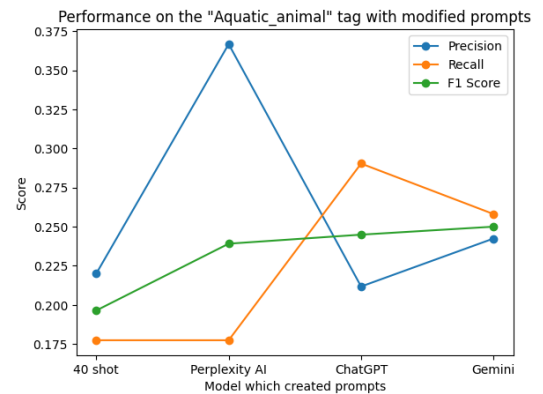
structures: (1) examples with explanations, (2) fill in the blank, and (3) error-augmented prompts. These three approaches were generated using Perplexity AI, ChatGPT, and Gemini respectively. I concluded that the 40-shot approach worked the best because it had the highest F1 score, but also acknowledged that each run type performed better than the others at identifying at least one token type. With further experimentation, I would have liked to see how a combination of prompt styles changed the performance of the LLM.

Section 3.4 Limitations As previously noted, resource limitations in the form of available credits significantly impacted the scope of the experiments. My original plan included the following runs: {0, 1, 5, 10, 20, 40, 100,

200} shots with original prompt structure, {0, 1, 5, 10, 20, 40, 100, 200} shots with a {fill in the blank, examples with explanations, error-augmented prompts} prompt structure, and an undetermined number of runs combining different numbers of shots and multiple prompt techniques. Unfortunately, due to credit constraints, executing the full range of experiments was not feasible. I used the OpenAI backend on each of the runs I executed. Using the Cohere platform was a potential way to expand expanding the run list, but my concerns regarding potential bias introduced by utilizing different backends led me to decide that it was best to maintain consistency for the sake of result purity. Additionally, I could have paid for my own API key, but without knowing how much cost that could incur, I decided against this approach.

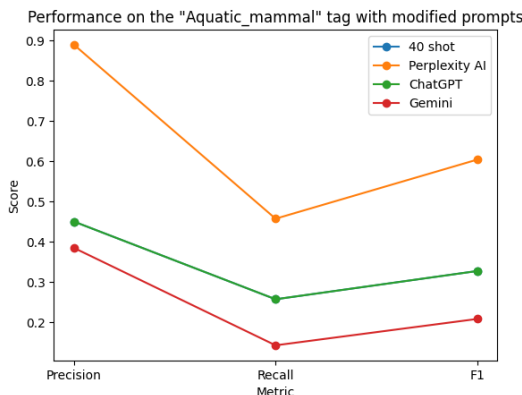


(a) Performance changes plotted by metric

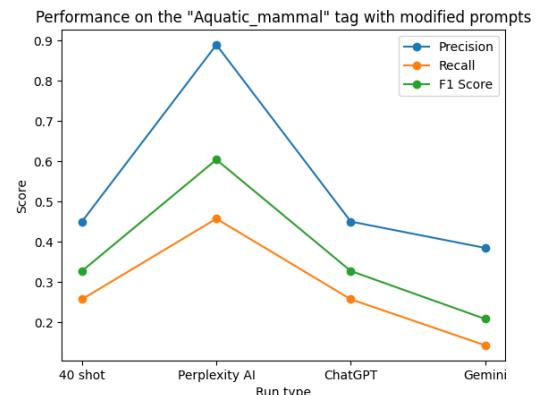


(b) Performance changes plotted by structure

Figure 13: Performance changes on the “Aquatic_animal” tag; generation code in Appendix Section A5

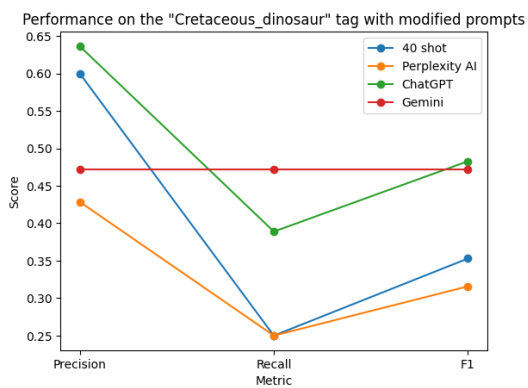


(a) Performance changes plotted by metric; NOTE: the 40 shot run and the ChatGPT run had the exact same metrics, which is why there are only 3 visible lines

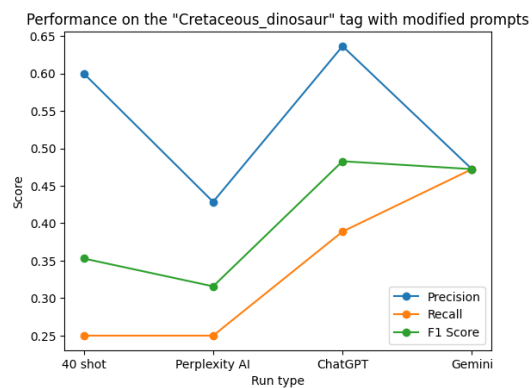


(b) Performance changes plotted by structure

Figure 14: Performance changes on the “Aquatic_mammal” tag; generation code in Appendix Section A5

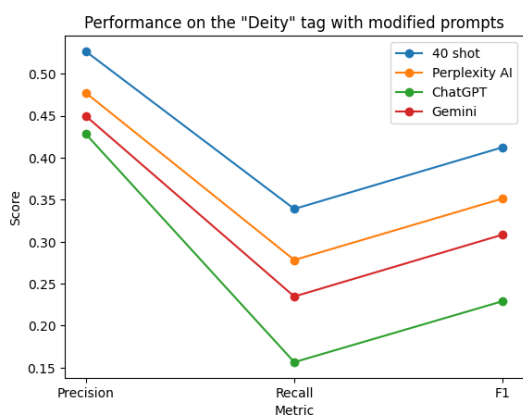


(a) Performance changes plotted by metric

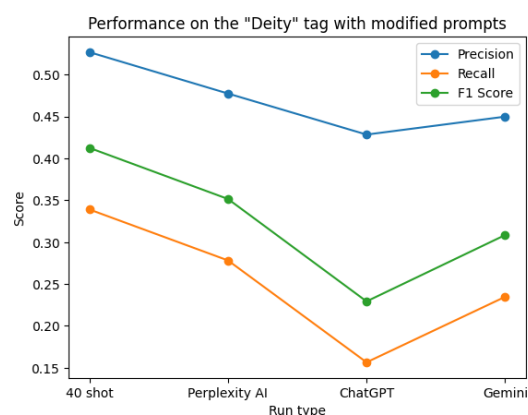


(b) Performance changes plotted by structure

Figure 15: Performance changes on the “Cretaceous_dinosaur” tag; generation code in Appendix Section A5

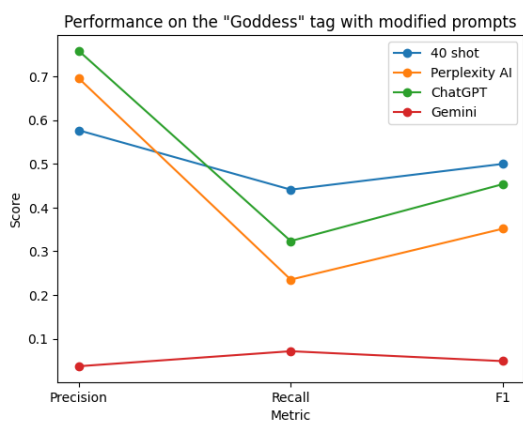


(a) Performance changes plotted by metric

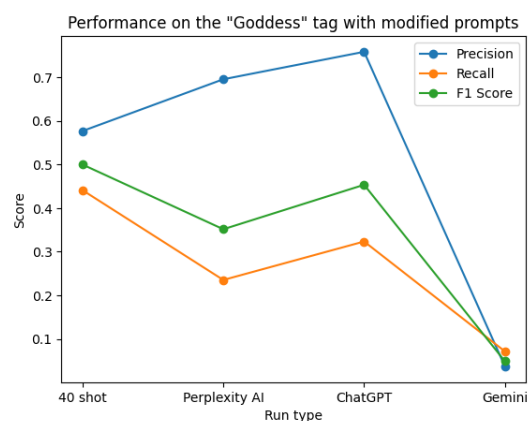


(b) Performance changes plotted by structure

Figure 16: Performance changes on the “Deity” tag; generation code in Appendix Section A5

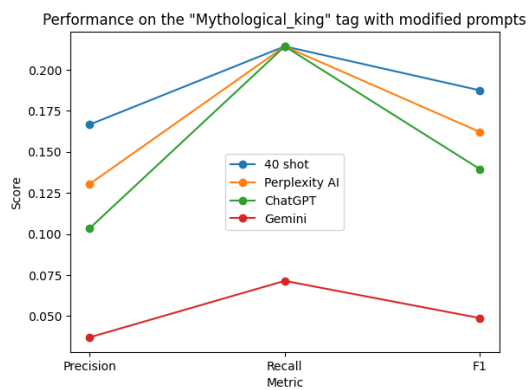


(a) Performance changes plotted by metric

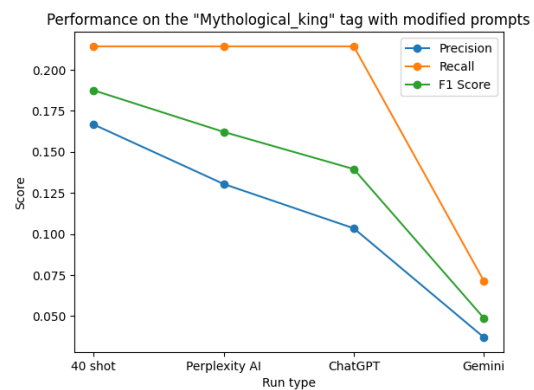


(b) Performance changes plotted by structure

Figure 17: Performance changes on the “Goddess” tag; generation code in Appendix Section A5



(a) Performance changes plotted by metric



(b) Performance changes plotted by structure

Figure 18: Performance changes on the “Mythological_king” tag; generation code in Appendix Section A5

Appendix: Code

Section A0: Original prompt structure code

```
1 # The baseline prompt structure format
2 messages = [
3     {'role': SYSTEM_STR, MSG_STR:
4         """"You will be given input text containing different types of entities that you will
5         label.
6         This is the list of entity types to label: Deity, Mythological_king,
7         Cretaceous_dinosaur, Aquatic_mammal, Aquatic_animal, Goddess.
8         Label the entities by surrounding them with tags like '<Cretaceous_dinosaur>
9         Beipiaognathus </Cretaceous_dinosaur>'."""
10    },
11    {'role': USER_STR, MSG_STR: """"Text: Once paired in later myths with her Titan brother
12    Hyperion as her husband, mild-eyed Euryphaessa, the far-shining one of the Homeric Hymn
13    to Helios, was said to be the mother of Helios (the Sun), Selene (the Moon), and Eos (
14    the Dawn)."""},
15    {'role': SYSTEM_STR, MSG_STR: """"Labels: Once paired in later myths with her Titan
16    brother <Deity> Hyperion </Deity> as her husband, mild-eyed Euryphaessa, the far-shining
17    one of the Homeric Hymn to Helios, was said to be the mother of Helios (the Sun), <
18    Goddess> Selene </Goddess> (the Moon), and <Goddess> Eos </Goddess> (the Dawn)."""},
19    {'role': USER_STR, MSG_STR: """"Text: From her ideological conception, Taweret was
20    closely grouped with (and is often indistinguishable from) several other protective
21    hippopotamus goddesses: Ipet, Reret, and Hedjet.\nLabels: """}
22 ]
23
24 # Returns the content of the example
25 def get_message(example):
26     return """"{}""".format(example['content'])
27
28 # A function which takes an example from the dataset as input, and returns a string that has
29 # tagged the text with labels in the given HTML-style format.
30 def convert_bio_to_prompt(example):
31     message = ""
32     current_label = None
33     # All we need is the token and its corresponding ner_string, we can ignore the rest of
34     # the information stored for each example
35     for token, ner_string in zip(example['tokens'], example['ner_strings']):
36         if ner_string != "0":
37             # Extract the label from the token
38             label = ner_string.split('-')[-1]
39             if label != current_label:
40                 if current_label:
41                     # Finishing the tagged entity, so add `</` to indicate the closing tag
42                     message += "</{}> ".format(current_label)
43                 # Otherwise, we are at the beginning of the tagged entity
44                 message += "<{}> ".format(label)
45                 current_label = label
46             message += token.split('-')[0] + " "
47     # If the current token is not a named entity but there was just one
48     elif current_label:
49         message += "</{}> ".format(current_label)
50         current_label = None
51     else:
52         message += token.split('-')[0] + " "
53     # If the example ends with a named entity
54     if current_label:
55         message += "</{}> ".format(current_label)
56     return message
```

```

46 # A function that takes the number of shots, dataset, list of entity types, and
    convert_bio_to_prompt function, and returns the chat_history (a list of maps) structured
    as in the example.
47 def get_chat_history(shots, dataset, entity_types_list, convert_bio_to_prompt_fn):
48     chat_history = []
49
50     # System message
51     system_message = {
52         'role': SYSTEM_STR, MSG_STR: "You will be given input text containing different
    types of entities that you will label. This is the list of entity types to label: {}.
    Label the entities by surrounding them with tags like '<{}> Entity </{}>'.format(
53         ", ".join(entity_types_list), list(entity_types_list)[0], list(entity_types_list)
    )[0])
54     }
55     chat_history.append(system_message)
56
57     for i in range(shots):
58         example = dataset[i]
59         # Get the text from each sample in the dataset
60         user_message = {
61             'role': USER_STR, MSG_STR: ""Text: {}"".format(get_message(example))
62         }
63         chat_history.append(user_message)
64
65         # Add the labels to the sample
66         system_message = {
67             'role': SYSTEM_STR, MSG_STR: ""Labels: {}"".format(convert_bio_to_prompt_fn(
    example))
68         }
69         chat_history.append(system_message)
70
71     return chat_history
72
73 # Function to handle punctuation which is not correctly attached to the closest token to the
    left
74 def move_punctuation_left(text):
75     # Define regular expression pattern to match punctuation followed by spaces
76     pattern = re.compile(r'\s*([.,;!?])\s*')
77     matches = re.finditer(pattern, text)
78     modified_text = ''
79     # Position of last match
80     last_match_end = 0
81     for match in matches:
82         punctuation = match.group(1)
83         # Get the start and end index of the match
84         start, end = match.span()
85         # Move the punctuation mark to touch the closest word on the left
86         modified_text += text[last_match_end:start].rstrip() + punctuation
87         # Update the position of the last match
88         last_match_end = end
89     # Append the remaining text after the last match
90     modified_text += text[last_match_end:]
91     return modified_text
92
93 def convert_response_to_bio(response):
94     labels = []
95     # bool: are we inside of a tag, bool: are we at the first position within the tag, bool:
    was the previous position in the tag, str: the name of the tag)
96     (inside, first, prev, tag) = (False, False, False, None)
97     for word in response.split():
98         # Ignore the starting labels

```

```

99         if "Labels:" in word:
100             continue
101         # At the end of a tag
102         if "</" in word:
103             rest = re.sub(r'<.+?>', '', word)
104             tag = (word.split('</')[1].split('>')[0]
105             if (word[-1] != '>' and any(char in string.punctuation for char in word.split('>
106             ')) and word.split('<')[0] == ''):
107                 # Reset everything because we are now outside of the tag
108                 (inside, first, prev, tag) = (False, False, False, None)
109                 continue;
110             # Inside of the tag
111             elif (rest != '' and prev):
112                 labels.append("I-" + tag)
113             # At the beginning
114             elif (rest != '' and not prev):
115                 labels.append("B-" + tag)
116                 (inside, first, prev, tag) = (False, False, False, None)
117                 continue
118             # At the beginning of a tag
119             elif "<" in word:
120                 (inside, first, prev, tag) = (True, True, False, (word.split('<')[1].split('>')
121                 [0])
122                 rest = re.sub(r'<.+?>', '', word)
123                 # This tag has multiple tokens that it is applied to
124                 if (rest != ''):
125                     labels.append("B-" + tag)
126                     first = False
127                     prev = True
128                 continue
129             if inside:
130                 if first:
131                     labels.append("B-" + tag)
132                     first = False
133                     prev = True
134                 else:
135                     labels.append("I-" + tag)
136             else:
137                 labels.append("O")
138
139         # Ensure that lonely punctuation marks are handled properly
140         tokens_pattern = re.compile(r'<[^\>]+>')
141         no_tags = re.sub(tokens_pattern, '', response)
142         no_punct_spaces = re.sub(r'\s(?:[?!."](?:\s|$))', r'\1', no_tags)
143         no_punct_spaces2 = move_punctuation_left(no_tags)
144
145         tokens = re.sub(r'([{}])'.format(r'\s*.,;!?\\s*'), r'\1 ', no_punct_spaces2).split()
146         if "Label" in tokens[0]:
147             tokens = tokens[1:]
148
149         return labels, tokens

```

Section A1: Perplexity AI “examples with explanations” prompt structure

The full output created by Perplexity AI as an example input for the “examples with explanations” prompt structure. See **Figure 2** and **Figure 3** in **Section 1** for the prompts given to ChatGPT to create this prompt structure.

```

1  messages = [
2  {'role': SYSTEM_STR, MSG_STR:
3  """You will be given input text containing different types of entities that you will
   label. This is the list of entity types to label: Deity, Mythological\_king, Cretaceous\
   \_dinosaur, Aquatic\_mammal, Aquatic\_animal, Goddess. Label the entities by surrounding

```

```

them with tags like '<Cretaceous\_dinosaur> Beipiaognathus </Cretaceous\_dinosaur>'.
Here are some examples with explanations:"""
},
{'role': USER_STR, MSG_STR:
"""Example 1:
Input text: 'Poseidon, the god of the sea, ruled over the oceans and all creatures that
dwelled within them, including dolphins and whales.'}
{'role': SYSTEM_STR, MSG_STR:
Labels: "<Deity> Poseidon </Deity>, the god of the sea, ruled over the oceans and all
creatures that dwelled within them, including <Aquatic\_mammal> dolphins </Aquatic\_
\_mammal> and <Aquatic\_mammal> whales </Aquatic\_mammal>."
Explanation: Poseidon is a deity, dolphins and whales are aquatic mammals."""
},
{'role': USER_STR, MSG_STR:
"""Example 2:
Input text: 'The infraclassis Carinacea includes most living species of regular sea
urchin, and fossil forms going back as far as the Triassic.'}
{'role': SYSTEM_STR, MSG_STR:
Labels: 'The infraclassis Carinacea includes most living species of regular <Aquatic\_
\_animal> sea urchin </Aquatic\_animal>, and fossil forms going back as far as the
Triassic.'
Explanation: A sea urchin is an aquatic animal."""
},
{'role': USER_STR, MSG_STR:
Now, try labeling the following input text:
Input text: "From her ideological conception, Taweret was closely grouped with (and is
often indistinguishable from) several other protective hippopotamus goddesses: Ipet,
Reret, and Hedjet."}

```

The change from the baseline code to generate this structure of prompt is in `get_chat_history`. This version extracts an explanation from the labeling and appends it to the system message.

```

def add_explanation():
    # Use regular expression to find content within <tag> tags
    match = re.search(r'<(.*?)>(.*?)</\1>', text)
    if match:
        # Extract tag name and content within the tags
        tag_name = match.group(1)
        content = match.group(2)
        # Create the explanation
        return f"Explanation: A {content.lower()} is a {tag_name}."
    else:
        return "No explanation found."

def get_chat_history(shots, dataset, entity_types_list, convert_bio_to_prompt_fn):
    # System message
    system_message = {
        'role': SYSTEM_STR, MSG_STR: """You will be given input text containing different
types of entities that you will label. This is the list of entity types to label: {}.
Label the entities by surrounding them with tags like '<{}> Beipiaognathus </{}>'. Here
are some examples with explanations:"""
    }.join(entity_types_list), list(entity_types_list)[0], list(entity_types_list)[0])
    chat_history.append(system_message)

    for i in range(shots):
        example = dataset[i]
        user_message = {
            'role': USER_STR, MSG_STR: """Text: {}""".format(get_message(example))
        }
        chat_history.append(user_message)

```



```

26         labeled = convert_bio_to_prompt_fn(example)
27         system_message = {
28             'role': SYSTEM_STR, MSG_STR: ""Labels: {} {}".format(labeled, add_explanation
29             (labeled))
30         }
31         chat_history.append(system_message)
32
33     return chat_history

```

Section A2: ChatGPT “fill in the blank” prompt structure

This is the example prompt generated by ChatGPT using the “fill in the blank” prompt structure. See **Figure 2** and **Figure 3** in **Section 1** for the prompts given to ChatGPT to create this prompt structure.

```

1 messages = [
2     {'role': SYSTEM_STR, MSG_STR:
3         ""Identify and label the entities in the following text. Replace the blank space with
4         the appropriate tags. This is the list of entity types to label: Deity,
5         Mythological_king, Cretaceous_dinosaur, Aquatic_mammal, Aquatic_animal, Goddess. Label
6         the entities by replacing the blank space with tags like '<Deity> Hyperion </Deity>'.""
7     },
8     {'role': USER_STR, MSG_STR: ""Text: Once paired in later myths with her Titan brother
9     Hyperion as her husband, mild-eyed Euryphaessa, the far-shining one of the Homeric Hymn
10    to Helios, was said to be the mother of Helios (the Sun), Selene (the Moon), and Eos (
11    the Dawn).""},
12    {'role': SYSTEM_STR, MSG_STR: ""Labels: Once paired in later myths with her Titan
13    brother <__> Hyperion <__> as her husband, mild-eyed Euryphaessa, the far-shining one of
14    the Homeric Hymn to Helios, was said to be the mother of Helios (the Sun), <__> Selene
15    <__> (the Moon), and <__> Eos <__> (the Dawn).""}]
16
17 message = f""Text: From her ideological conception, Taweret was closely grouped with (and
18 is often indistinguishable from) several other protective hippopotamus goddesses: Ipet,
19 Reret, and Hedjet.
20 Labels: ""

```

The most important change to the baseline code to create the ChatGPT fill-in-the-blank style prompt is to the function `convert_bio_to_prompt`. Instead of adding the proper entity tags, we add “_” everywhere.

```

1 # The modified convert_bio_to_prompt code to generate prompts of the same structure as the
2 # example created by ChatGPT
3 def convert_bio_to_prompt(example):
4     message = ""
5     # The labels will all be replaced with "_" to create the fill-in-the-blank structure
6     current_label = "_"
7     for token, ner_string in zip(example['tokens'], example['ner_strings']):
8         if ner_string != "0":
9             label = ner_string.split('-')[-1] # Extracting the entity label from the token
10            if label != current_label:
11                if current_label:
12                    message += "<{}> ".format(current_label)
13                    # Add "_" instead of the actual tag to create fill-in-the-blank prompt
14                    message += "<{}> ".format("_")
15                    current_label = "_"
16                message += token.split('-')[0] + " "
17            # If the current token is not a non-"0" tag but there was just one
18            elif current_label:
19                # Once again add "_" instead of actual tag to create fill-in-the-blank prompt
20                message += "<{}> ".format("_")
21                current_label = None

```

```

21         else:
22             message += token.split('-')[0] + " "
23
24     # Check if the ending is a non-"0" tag
25     if current_label:
26         message += "<{}>".format(current_label)
27     return message

```

Section A3: Gemini “error-augmented prompts” prompt structure

This is the example prompt generated by Gemini using the ‘error-augmented prompts’ prompt structure. See **Figure 2** and **Figure 3** in **Section 1** for the prompts given to ChatGPT to create this prompt structure.

```

1     original_text = "Once paired in later myths with her Titan brother Hyperion as her
    husband, mild-eyed Euryphaessa, the far-shining one of the Homeric Hymn to Helios, was
    said to be the mother of Helios (the Sun), Selene (the Moon), and Eos (the Dawn)."
2     original_label = "Once paired in later myths with her Titan brother <Deity> Hyperion </
    Deity> as her husband, mild-eyed Euryphaessa, the far-shining one of the Homeric Hymn to
    Helios, was said to be the mother of Helios (the Sun), <Goddess> Selene </Goddess> (the
    Moon), and <Goddess> Eos </Goddess> (the Dawn)."
3
4     error_text_1 = original_text.replace("Hyperion", "Hipyrrion").replace("Selene", "sElenE")
5     err_label_1 = original_label.replace("Hyperion", "Hipyrrion").replace("Selene", "sElenE")
6
7     clean_prompt = {'role': USER_STR, MSG_STR: f"Text: {original_text}"}
8     clean_response = {'role': SYSTEM_STR, MSG_STR: f"Labels: {original_label}"}
9
10    error_prompt = {'role': USER_STR, MSG_STR: f"Text: {error_text_1}"}
11    error_response = {'role': SYSTEM_STR, MSG_STR: f"Labels: {err_label_1}"}
12
13    messages = [
14        {'role': SYSTEM_STR, MSG_STR:
15            """You will be given input text containing different types of entities that you will
            label. This is the list of entity types to label: Deity, Mythological_king,
            Cretaceous_dinosaur, Aquatic_mammal, Aquatic_animal, Goddess. Label the entities by
            surrounding them with tags like '<Cretaceous_dinosaur> Beipiaognathus </
            Cretaceous_dinosaur>'."""
16        },
17        clean_prompt, clean_response, error_prompt, error_response
18    ]

```

The significant change to the baseline code to create the Gemini error-augmented prompt structure style is to `get_chat_history`. Instead of using a basic pair of input text and labeled text, each pair adds either spelling mistakes, strange capitalization, filler words, or other noise. This is intended to help the model learn to parse through the noise and still correctly label each token with the proper tag. All other functions are the same as in the baseline implementation.

```

1 # A package to generate random sentences: https://pypi.org/project/wonderwords/
2 from wonderwords import RandomSentence
3
4 # This function creates the chat history using the error-augmented prompts structure. The
    idea is to inject errors, misspellings, and filler words into the prompts in the hopes
    that the model will still learn to pick up the correct context in order to label each
    token correctly.
5 def get_chat_history(shots, dataset, entity_types_list, convert_bio_to_prompt_fn):
6     original_text = "Once paired in later myths with her Titan brother Hyperion as her
    husband, mild-eyed Euryphaessa, the far-shining one of the Homeric Hymn to Helios, was
    said to be the mother of Helios (the Sun), Selene (the Moon), and Eos (the Dawn)."
7     original_label = "Once paired in later myths with her Titan brother <Deity> Hyperion </
    Deity> as her husband, mild-eyed Euryphaessa, the far-shining one of the Homeric Hymn to

```

```
Helios, was said to be the mother of Helios (the Sun), <Goddess> Selene </Goddess> (the Moon), and <Goddess> Eos </Goddess> (the Dawn)."
```

```
# Error 1: Misspelling and Case Change. This message contains misspelled words and words with weird casing. The LLM should be able to handle these malformed messages in the same way as the basic ones.
```

```
error_text_1 = original_text.replace("Hyperion", "Hipyrrion").replace("Selene", "sELenE")  
err_label_1 = original_label.replace("Hyperion", "Hipyrrion").replace("Selene", "sELenE")
```

```
# Error 2: Word Substitution and Noise Words. This message substitutes some words and adds other irrelevant words, none of which should be tagged as a special entity.
```

```
error_text_2 = original_text.replace("Titan brother", "large sibling").replace("far-shining", "bright and beautiful") + " In addition, a random sea monster also appeared."  
error_label_2 = original_label.replace("Titan brother", "large sibling").replace("far-shining", "bright and beautiful") + " In addition, a random sea monster also appeared."
```

```
# Error 3: Irrelevant Entity. This message contains both noise words and an entity which doesn't fit under any of the tags we have specified.
```

```
error_text_3 = "The is some text which doesn't do anything. " + original_text + " Also, King Arthur was said to be a great leader."
```

```
error_label_3 = "The is some text which doesn't do anything. " + original_label + " Also, King Arthur was said to be a great leader."
```

```
# Create separate prompts with original and error-augmented data
```

```
clean_prompt = {'role': USER_STR, MSG_STR: f"Text: {original_text}"}  
clean_response = {'role': SYSTEM_STR, MSG_STR: f"Labels: {original_label}"}  
  
error_prompt_1 = {'role': USER_STR, MSG_STR: f"Text: {error_text_1}"}  
error_response_1 = {'role': SYSTEM_STR, MSG_STR: f"Labels: {err_label_1}"}  
  
error_prompt_2 = {'role': USER_STR, MSG_STR: f"Text: {error_text_2}"}  
error_response_2 = {'role': SYSTEM_STR, MSG_STR: f"Labels: {error_label_2}"}  
  
error_prompt_3 = {'role': USER_STR, MSG_STR: f"Text: {error_text_3}"}  
error_response_3 = {'role': SYSTEM_STR, MSG_STR: f"Labels: {error_label_3}"}  
  
# Join the starting message with instructions to the rest of the formed messages  
messages = [  
    {'role': SYSTEM_STR, MSG_STR:  
        """You will be given input text containing different types of entities that you will label. This is the list of entity types to label: Deity, Mythological_king, Cretaceous_dinosaur, Aquatic_mammal, Aquatic_animal, Goddess. Label the entities by surrounding them with tags like '<Cretaceous_dinosaur> Beipiaognathus </Cretaceous_dinosaur>'."""  
    },  
    clean_prompt, clean_response, error_prompt_1, error_response_1, error_prompt_2, error_response_2, error_prompt_3, error_response_3]  
  
s = RandomSentence()  
  
# For the rest of the prompts, we want to add filler words and sentences, hoping that the model will learn to handle/ignore these irrelevant sections without compromising its ability to identify the tokens which need to be specially labelled.  
for i in range(shots):  
    example = dataset[i]  
    # Create a random (but grammatically correct) sentence  
    sent = s.sentence()  
    # To introduce a bit of extra pseudo-randomness, alternate if the filler sentence is at the beginning or the end of the proper message  
    if i % 2 == 0:  
        user_message = {
```

```

51         'role': USER_STR, MSG_STR: ""Text: {} {}".format(sent, get_message(
example))
52     }
53     system_message = {
54         'role': SYSTEM_STR, MSG_STR: ""Labels: {} {}".format(sent,
convert_bio_to_prompt_fn(example))
55     }
56     else:
57         user_message = {
58             'role': USER_STR, MSG_STR: ""Text: {} {}".format(get_message(example),
sent)
59         }
60         system_message = {
61             'role': SYSTEM_STR, MSG_STR: ""Labels: {} {}".format(
convert_bio_to_prompt_fn(example), sent)
62         }
63         # Append all of the messages as normal
64         messages.append(user_message)
65         messages.append(system_message)
66
67     return messages

```

Section A4: Experimental results code

```

1 # The eval code is the same as the baseline implementation
2 def run_eval(dataset, shots, backend):
3
4     for example in tqdm(dataset, total=len(dataset), desc="Evaluating", position=tqdm.
_get_free_pos()):
5
6         # String list of labels (BIO)
7         true_labels = [labels_int2str[l] for l in example['ner_tags']]
8         example_tokens = example['tokens']
9
10        response_text = call_api_openai(shots, example) if backend == "openai" else
call_api_cohere(shots, example)
11
12        # String list of predicted labels (BIO)
13        predictions, generated_tokens = convert_response_to_bio(response_text)
14
15        if len(example['ner_strings']) != len(predictions):
16            predictions = example['ner_strings']
17
18        # Handle case where the generated text doesn't align with the input text.
19        # Basically, we'll eval everything up to where the two strings start to diverge.
20        # We relax this slightly by ignoring punctuation (sometimes we lose a paren or
something,
21        # but that's not catastrophic for eval/tokenization).
22        # Just predict 'O' for anything following mismatch.
23        matching_elements = [strip_punct(i) == strip_punct(j) for i, j in zip(example_tokens,
generated_tokens)]
24
25        if False in matching_elements:
26            last_matching_idx = matching_elements.index(False)
27        else:
28            last_matching_idx = min(len(generated_tokens), len(example_tokens))
29
30        predictions = predictions[:last_matching_idx] + ['O']*(len(example_tokens)-
last_matching_idx)
31        metric.add(predictions=predictions, references=true_labels)
32
33    return metric.compute()

```

```

34 # Run the eval on the entire dev set
35 dev_examples_to_take = 0
36
37 dev_set = data_splits['dev']
38 if dev_examples_to_take > 0:
39     dev_set = data_splits['dev'].select(range(dev_examples_to_take))
40
41 # The number of shots was changed according to the specific run
42 for num_shots in [5]:
43     print(f"shots: {num_shots}")
44     result = run_eval(dev_set, shots=num_shots, backend='openai')
45     print(result)
46
47 #####
48
49 # The rest of the code is written to create the plots for the experimental results
50 # The following two sections of code create the graphs to display the performance changes as
    a result of changing the prompt structure.
51 # These numbers were extracted directly from the output after the evaluation was run. This
    plot and the plot below display the same data, but with different x-axis values. This
    plot has the performance metrics on the x-axis.
52 labels = ["Accuracy", "Precision", "Recall", "F1 Score"]
53 perplexity = [0.964, 0.363, 0.242, 0.298]
54 gpt = [0.956, 0.370, 0.256, 0.302]
55 gemini = [0.958, 0.347, 0.263, 0.290]
56
57 plt.plot(labels, perplexity, "-o", label="Perplexity AI")
58 plt.plot(labels, gpt, "-o", label="ChatGPT")
59 plt.plot(labels, gemini, "-o", label="Gemini")
60 plt.legend(loc="best")
61 plt.title('Performance comparision between Perplexity AI, Gemini, and GPT prompts')
62 plt.xlabel('Performance metric')
63 plt.ylabel('Score')
64 plt.show()
65
66 #####
67
68 # These numbers were extracted directly from the output after the evaluation was run. This
    plot has the run type on the x-axis.
69 acc = [0.964, 0.956, 0.958]
70 prec = [0.363, 0.370, 0.347]
71 rec = [0.242, 0.256, 0.263]
72 f = [0.298, 0.302, 0.290]
73 labels = ["Perplexity AI", "ChatGPT", "Gemini"]
74
75 plt.plot(labels, acc, "-o", label="Accuracy")
76 plt.plot(labels, prec, "-o", label="Precision")
77 plt.plot(labels, rec, "-o", label="Recall")
78 plt.plot(labels, f, "-o", label="F1")
79 plt.legend(loc="best")
80 plt.title('Performance comparision between Perplexity AI, Gemini, and GPT prompts')
81 plt.xlabel('Model which created prompts')
82 plt.ylabel('Score')
83 plt.show()
84
85 #####
86
87 # The following code creates the graphs to display the performance changes as a result of
    increasing the number of shots
88 # These numbers were extracted directly from the output after the evaluation was run. This
    plot and the plot below display the same data, but with different x-axis values. This

```

```

    plot has the performance metrics on the x-axis.
89 zerosht = [0.965, 0.460, 0.239, 0.296]
90 onesht = [0.959, 0.322, 0.300, 0.311]
91 fivesht = [0.965, 0.389, 0.239, 0.296]
92 tensht = [0.965, 0.329, 0.212, 0.258]
93 twentysht = [0.948, 0.366, 0.258, 0.302]
94 fortysht = [0.946, 0.441, 0.306, 0.361]
95 labels = ["Accuracy", "Precision", "Recall", "F1 Score"]
96
97 plt.plot(labels, zerosht, "-o", label='0 Shot')
98 plt.plot(labels, onesht, "-o", label='1 Shot')
99 plt.plot(labels, fivesht, "-o", label='5 Shot')
100 plt.plot(labels, tensht, "-o", label='10 Shot')
101 plt.plot(labels, twentysht, "-o", label='20 Shot')
102 plt.plot(labels, fortysht, "-o", label='40 Shot')
103 plt.legend(loc="best")
104 plt.title('Performance changes as a function of number of few-shot examples')
105 plt.xlabel('Performance metric')
106 plt.ylabel('Score')
107 # for e in [onesht, fivesht, tensht, twentysht, fortysht]:
108 #     for x, y in zip(labels, e):
109 #         plt.annotate(text=str(y), xy=(x, y))
110 plt.show()
111
112 #####
113
114 # These numbers were extracted directly from the output after the evaluation was run. This
    plot has the number of shots on the x-axis.
115 accuracy = [0.965, 0.959, 0.965, 0.965, 0.948, 0.946]
116 precision = [0.460, 0.322, 0.389, 0.329, 0.366, 0.441]
117 recall = [0.239, 0.300, 0.239, 0.212, 0.258, 0.306]
118 f1 = [0.296, 0.311, 0.296, 0.258, 0.302, 0.361]
119 labels = ["0 shot", "1 shot", "5 shots", "10 shots", "20 shots", "40 shots"]
120
121 plt.plot(labels, accuracy, "-o", label="Accuracy")
122 plt.plot(labels, precision, "-o", label="Precision")
123 plt.plot(labels, recall, "-o", label="Recall")
124 plt.plot(labels, f1, "-o", label="F1 Score")
125 plt.legend(loc="best")
126 plt.title('Performance changes as a function of number of few-shot examples')
127 plt.xlabel('Number of shots')
128 plt.ylabel('Score')
129 plt.show()

```

Section A5: Analysis and Discussion code

Code used to generate Figures 10 - 12

```

1 # Code used to create the graphs of tag distribution across the dev and training sets. This
    code was used in the analysis of how the distribution of tags in the datasets affected
    the performance of different prompting styles
2 import json
3 import re
4 import matplotlib.pyplot as plt
5 from collections import Counter
6
7 tags = Counter()
8 tags_train = Counter()
9 tags_dev = Counter()
10
11 # Go through the dev and training datasets and collect the ner strings for each sample
12 for entry in data_splits['dev']:

```



```

13 tags.update(entry['ner_strings'])
14 tags_dev.update(entry['ner_strings'])
15 for entry in data_splits['train']:
16     tags.update(entry['ner_strings'])
17     tags_train.update(entry['ner_strings'])
18
19 total = sum(tags.values())
20
21 # Sort by the name of the tag, ignoring the B- and I-
22 def custom_sort(tag):
23     return tag[2:]
24
25 def plot_each(cnt, st):
26
27     # Extract all labels that aren't 'O'
28     non0 = Counter({k: c for k, c in cnt.items() if k != "O"})
29
30     # The x-axis is the tags sorted in alphabetical order
31     labels = list(non0.keys())
32     sorted_labels = sorted(labels, key=custom_sort)
33     sorted_counts = [non0[label] for label in sorted_labels]
34
35     plt.bar(sorted_labels, sorted_counts)
36     plt.xlabel('NER labels')
37     plt.xticks(rotation=90)
38     plt.ylabel('Occurrences')
39     if st == "all":
40         plt.title('Occurrences of all NER labels without "O" across all datasets')
41     else:
42         plt.title('Occurrences of all NER labels without "O" in the ' + st + ' set')
43
44     plt.show()
45
46 plot_each(tags_train, "training")
47 plot_each(tags_dev, "dev")
48 plot_each(tags, "all")

```

Code used to create graphs seen in Figures 13 - 18

```

1 # A function to create the graphs for each tag. These graphs were used in the analysis of
  # performance of different prompting styles.
2 def plot_by_tag(tag, tag_prec, tag_rec, tag_f1, forty, perp, chatpt, gem):
3     # These labels are the types of runs, the first graph uses them on the x-axis. The y-
  # axis is the performance metric score scale.
4     labels = ["40 shot", "Perplexity AI", "ChatGPT", "Gemini"]
5     plt.plot(labels, tag_prec, "-o", label="Precision")
6     plt.plot(labels, tag_rec, "-o", label="Recall")
7     plt.plot(labels, tag_f1, "-o", label="F1 Score")
8     plt.legend(loc="best")
9     plt.title('Performance on the "' + tag + '" tag with modified prompts')
10    plt.xlabel('Run type')
11    plt.ylabel('Score')
12    plt.show()
13
14    # These labels are the performance metrics, the second graph uses them on the x-axis.
  # The y-axis is the performance metric scores.
15    labels = ["Precision", "Recall", "F1"]
16    plt.plot(labels, forty, "-o", label="40 shot")
17    plt.plot(labels, perp, "-o", label="Perplexity AI")
18    plt.plot(labels, chatpt, "-o", label="ChatGPT")
19    plt.plot(labels, gem, "-o", label="Gemini")
20    plt.legend(loc="best")

```

```

21 plt.title('Performance on the "' + tag + '" tag with modified prompts')
22 plt.xlabel('Metric')
23 plt.ylabel('Score')
24 plt.show()

```

References

- Bsharat S. M., Myrzakhan A., Shen Z., 2024, Principled Instructions Are All You Need for Questioning LLaMA-1/2, GPT-3.5/4 ([arXiv:2312.16171](#))
- Chen B., Zhang Z., Langrené N., Zhu S., 2023, Unleashing the potential of prompt engineering in Large Language Models: a comprehensive review ([arXiv:2310.14735](#))
- Eliot L., 2023, Latest Prompt Engineering Technique Inventively Transforms Imperfect Prompts into Superb Interactions for Using Generative AI, <https://www.forbes.com/sites/lanceeliot/2023/07/26/latest-prompt-engineering-technique-inventively-transforms-imperfect-prompts-into-superb-interactions-for-using-generative-ai/?sh=5af16bdb2c31>
- Liu P., Yuan W., Fu J., Jiang Z., Hayashi H., Neubig G., 2021, Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing ([arXiv:2107.13586](#))
- Reynolds L., McDonell K., 2021, Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm ([arXiv:2102.07350](#))
- Tang L., Peng Y., Wang Y., Ding Y., Durrett G., Rousseau J. F., 2023, Less Likely Brainstorming: Using Language Models to Generate Alternative Hypotheses ([arXiv:2305.19339](#))
- The New Stack 2023, Prompt Engineering: Get LLMs to Generate the Content You Want, <https://thenewstack.io/prompt-engineering-get-llms-to-generate-the-content-you-want/>
- Weird_Foundation5044 2023, Is it just me or is Prompt Engineering basically..., https://www.reddit.com/r/LocalLLaMA/comments/182pjyg/is_it_just_me_or_is_prompt_engineering_basically/
- Weng L., 2023, lilianweng.github.io