

Project Proposal: The Practicality of Prompt Engineering

Anonymous ACL submission

Abstract

This paper examines the practicality of prompt engineering in improving the performance of Large Language Models (LLMs). Through empirical analysis, we evaluate the trade-offs between costs and benefits of prompting using novel metrics. Different prompting methods are assessed using standardized tasks and both modern and older models.

Introduction

Prompt engineering, the practice of developing specialized prompts and queries to improve the accuracy of Large Language Models after training, is a prominent topic of interest in the NLP community, and among the general public. The practice is believed to allow for improvements in LLM performance a variety of domains without investment in underlying training (Martineau, 2021). It is not, however, without its critics. Some commentators believe that the practice will become irrelevant as models grow larger and more powerful, becoming more capable of directly interpreting a user's intent. (Ethan Mollick [@emollick], 2023). Others question the need for specialized professionals or training to attain minimal improvements which are often not repeatable across domain types and contexts (Shackell, 2023; Acar, 2023).

Despite such controversy, it is difficult to find empirical and quantitative analyses of the trade-off between costs and accuracy benefits associated with advanced prompting. Papers introducing new prompting techniques often only include performance benchmarks concerning the techniques' efficacy, typically within a limited question domain. Some authors briefly mention problems associated with human-tailored problems, such as the increased complexity induced by prompt-chaining and limitations on creativity and randomness (Wu et al., 2022), and others suggest the automation of prompting (Diao et al., 2023). Online marketplaces

such as Promptbase (noa) provide input token costs for prompt texts sold on the platform, but do not provide any other information.

This paper uses several metrics to evaluate the costs of prompt engineering methods systematically, and analyzes the tradeoffs inherent in their application to standardized data. Such an assessment is valuable on several dimensions. Beyond quantifiably testing the practicality of prompt engineering as a whole, it can be used to compare the performance of different approaches, useful in a world where so many competing techniques are available. I also offer a new look at these methods in a period long after ideas were introduced and in an environment with greater capabilities in underlying models. Finally, I survey and introduce some useful measures of costs and complexity such as the ratio of interaction length with prompting to the length of an accepted human-generated answer.

Metrics

Accuracy

Improved accuracy can be a benefit of prompt engineering. I report:

- Correct/Incorrect accuracy at the point a technique has been fully implemented (the end of the chain of thought, etc., modelling the real world)

Length

The length of responses and interactions can effect the practicality of prompting. It could indicate that a model is carefully and correctly solving through the steps of a problem. It can also impose time and financial costs to users, or become an indicator of degraded performance as models sometimes tend to go off on tangents or repeat themselves (to the extent some platforms have imposed length limitations) (Mann). I report:

- Length of the entire interaction in tokens

078	• Financial cost of the entire interaction in to-	accuracy gains for each method. Higher accuracy	124
079	kens	methods are likely to be more interesting.	125
080	• Length of the entire interaction in tokens rel-	Below items are listed in order of increasing	126
081	ative to the length of the baseline task + a	complexity/human intervention:	127
082	human/solved out/generally accepted as cor-	• Zero-Shot Control Baseline: Providing the	128
083	rect answer. How much is the engineering	question/task directly. (Sometimes called "Di-	129
084	stretching the interaction out? Is this stretch-	rect Prompting")	130
085	ing adding value/improving accuracy?	• Zero-Shot Chain of Thought Prompting: Ex-	131
086	• Length of the entire interaction in time (sec-	isting literature mentions several examples of	132
087	onds). This can include time writing a re-	this: a simple one to append to every initial	133
088	sponse, waiting for a response, or reviewing	question/task, found to be optimal through au-	134
089	a response - this level of granular data may	tomated testing is "Let's work this out in a	135
090	be hard to collect, but it might be possible	step by step way to be sure we have the right	136
091	to look at human assessments of time spent	answer." (Hebenstreit et al., 2023; Zhou et al.,	137
092	on these activities. Attempts will be made to	2022)	138
093	have queries to models made during off-peak	• Tree of Thought Prompting: This prepends the	139
094	hours to minimize confounding due to server	following to the task/question: "Imagine three	140
095	load, connectivity issues, etc.	different experts are answering this question.	141
096	Complexity	All experts will write down 1 step of their	142
097	Similarly to length, complexity could be an indica-	thinking, then share it with the group. Then	143
098	tor of high accuracy. However, it has substantial	all experts will go on to the next step, etc. If	144
099	costs in potentially making review more difficult. I	any expert realises they're wrong at any point	145
100	report:	then they leave. The question is..." (Hulbert,	146
101	• Vocabulary complexity, sentence length, read-	2023)	147
102	ability scores of responses	• Chain of Verification Prompting: The LLM	148
103	• Ratio or difference of these scores in prompts	is prompted a series of times to produce an	149
104	vs. responses, responses vs. accepted/outside	initial baseline response, write its own veri-	150
105	correct answer	fication questions, answer those verification	151
106	• Human assessment of need for specialized	questions, and write a final, verified response.	152
107	knowledge/difficulty of implementation of	(Dhuliawala et al., 2023)	153
108	the technique. This could be task specific	• Few-Shot Prompting: The prompter provides	154
109	(done for some novel real world example ques-	a few examples of successfully/answered	155
110	tions/prompting scenarios), or it could be done	questions or tasks before the main ques-	156
111	overall based on a pre-existing description of	tion/task. Despite the work involved in im-	157
112	the technique. Perhaps a balance of both is	plementing this method, I believe it has po-	158
113	best.	tential to be effective for two reasons. First,	159
114	• Human assessment of output complexity (ease	prior evidence indicates that larger and more	160
115	of evaluating results). This could be task spe-	modern language models benefit more from	161
116	cific (done for some novel real world example	few-shot learning, potentially making this a	162
117	questions/prompting scenarios), or it could be	consistently useful technique. (Brown et al.,	163
118	done overall based on a pre-existing examples	2020) Second, earlier research has found that	164
119	of the technique. Perhaps a balance of both is	the formatting and input/label space distribu-	165
120	best.	tion is more important than example correct-	166
121	Prompting Methods Assessed	ness, meaning this method is somewhat robust	167
122	This list is subject to change, pending further as-	to human error. (Min et al., 2022)	168
123	essment of implementation difficulty and potential	Data	169
		To evaluate the LLMs, I attempt to use tasks that	170
		are both general-purpose and close to practical,	171

real-world applications. In this spirit, I use GRE General Test questions (likely from a purchased prep book/or practice exams recently published online to minimize contamination), as well as HumanEval coding problems. For multiple choice questions, I add formatting to the prompt to indicate what should be done to represent a final answer "When you are ready, please put your final, letter answer in the form of a single capital letter, enclosed in parentheses. For example, if you think the answer is A, please write (A)." HumanEval uses the pass at k metric to automatically assess the probability a solution is correct given a set of unit tests. (Chen et al., 2021)

I perform the analysis on one cutting-edge model, to get the current state of things, and one older model, closer to the time that these techniques were introduced. This will allow one picture of the changing costs and benefits of advanced prompting, a trend that may even be extrapolated into the future if current LLM scaling laws continue to hold. As the most widely used models and the ones behind much original work in the field, I select two models from the OpenAI series: GPT-4 and text-davinci-003 on OpenAI playground. Should text-davinci-003 (a legacy model) become unavailable during the course of the project, or should I encounter other difficulties, I will use the simplest/smallest model available, likely GPT-3.5 (something to note is that models older than GPT-3.5 have, in the past, scored 0% for accuracy on coding problems - I will need to test all of my evaluations quickly to get a sense of their feasibility before scaling up). All of these models are available both via the OpenAI API and in web interfaces - where possible I will use web interfaces to limit resources required.

To the extent possible, I will report accuracy scores on the domain dataset as they are in the original paper. It may also be possible to use any prompts, LLM responses, and correct responses provided along with original papers to calculate other simple metrics such as response lengths and complexity. However, some metrics (time taken, etc.) will require my own evaluations.

Analyses

I provide summary statistics of the metrics for each prompting method by model. In cases where human/textual assessment and comments have been provided, it might be interesting to use NLP methods to evaluate responses (ex: for sentiment).

Responsibilities

I am the sole author of this proposal, but I am open to working with others with similar interests. Expanding the group would be helpful in order to improve the project.

Limitations

It was difficult to select prompt engineering methods to try for this paper, and there is potential for my choice of methods to be somewhat biased. I mostly picked methods based my perception of their popularity and on their ease of implementation. If anything, this may lead to an underestimation of costs.

Just as my evaluation comes at a time with significantly more capable LLMs relative to that where much work began on prompting, I expect the underlying calculus concerning prompting to continue to change in the future. However, I again only expect relative costs of complex engineering to increase as models get better.

Another potential problem is the extent that prompting techniques have been absorbed into default LLM behavior, likely through reinforcement learning. GPT-4 in particular does seem to implement chain-of-thought methods when prompted with a sufficiently complex problem. In this environment, this paper become less of an evaluation of prompting techniques themselves, but more of an evaluation of their manual implementation.

Finally, though I have taken steps to limit it, data contamination remains a real concern. The evaluations I use may have been used to train the LLMs, or the questions/tasks from them may have been introduced through reinforcement learning based on other evaluations. On the other hand, this seems unlikely to bias the results for any one particular prompting method relative to the others - comparisons between them are still likely to be useful.

Acknowledgements

The template for this document was adapted by Jordan Boyd-Graber, Naoaki Okazaki, and Anna Rogers.

References

PromptBase.

267	Oguz A. Acar. 2023. AI Prompt Engineering Isn't the	Dave Hulbert. 2023. Using Tree-of-Thought Prompting	324
268	Future . <i>Harvard Business Review</i> . Section: Technol-	to boost ChatGPT's reasoning . Original-date: 2023-	325
269	ogy and analytics .	05-22T19:03:27Z.	326
270	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	Jyoti Mann. Microsoft limits Bing chat exchanges and	327
271	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	conversation lengths after 'creepy' interactions with	328
272	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	some users .	329
273	Askell, Sandhini Agarwal, Ariel Herbert-Voss,	Kim Martineau. 2021. What is prompt tuning?	330
274	Gretchen Krueger, Tom Henighan, Rewon Child,		
275	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe,	331
276	Clemens Winter, Christopher Hesse, Mark Chen, Eric	Mike Lewis, Hannaneh Hajishirzi, and Luke Zettle-	332
277	Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,	moyer. 2022. Rethinking the Role of Demonstra-	333
278	Jack Clark, Christopher Berner, Sam McCandlish,	tions: What Makes In-Context Learning Work?	334
279	Alec Radford, Ilya Sutskever, and Dario Amodei.	ArXiv:2202.12837 [cs] .	335
280	2020. Language Models are Few-Shot Learners .		
281	ArXiv:2005.14165 [cs] .	Cameron Shackell. 2023. Prompt engineering: is being	336
282	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming	an AI 'whisperer' the job of the future or a short-lived	337
283	Yuan, Henrique Ponde de Oliveira Pinto, Jared Ka-	fad?	338
284	plan, Harri Edwards, Yuri Burda, Nicholas Joseph,	Tongshuang Wu, Michael Terry, and Carrie Jun Cai.	339
285	Greg Brockman, Alex Ray, Raul Puri, Gretchen	2022. AI Chains: Transparent and Controllable	340
286	Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-	Human-AI Interaction by Chaining Large Language	341
287	try, Pamela Mishkin, Brooke Chan, Scott Gray,	Model Prompts . In <i>CHI Conference on Human Fac-</i>	342
288	Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz	tors in Computing Systems , pages 1–22, New Orleans	343
289	Kaiser, Mohammad Bavarian, Clemens Winter,	LA USA. ACM.	344
290	Philippe Tillet, Felipe Petroski Such, Dave Cum-		
291	mings, Matthias Plappert, Fotios Chantzis, Eliza-	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han,	345
292	beth Barnes, Ariel Herbert-Voss, William Hebgen	Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy	346
293	Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie	Ba. 2022. Large Language Models are Human-Level	347
294	Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,	Prompt Engineers .	348
295	William Saunders, Christopher Hesse, Andrew N.		
296	Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan		
297	Morikawa, Alec Radford, Matthew Knight, Miles		
298	Brundage, Mira Murati, Katie Mayer, Peter Welin-		
299	der, Bob McGrew, Dario Amodei, Sam McCandlish,		
300	Ilya Sutskever, and Wojciech Zaremba. 2021. Eval-		
301	uating Large Language Models Trained on Code .		
302	ArXiv:2107.03374 [cs] version: 2.		
303	Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu,		
304	Roberta Raileanu, Xian Li, Asli Celikyilmaz, and		
305	Jason Weston. 2023. Chain-of-Verification Re-		
306	duces Hallucination in Large Language Models .		
307	ArXiv:2309.11495 [cs] .		
308	Shizhe Diao, Pengcheng Wang, Yong Lin, and		
309	Tong Zhang. 2023. Active Prompting with		
310	Chain-of-Thought for Large Language Models .		
311	ArXiv:2302.12246 [cs] .		
312	Ethan Mollick [@emollick]. 2023. I have a strong sus-		
313	picion that "prompt engineering" is not going to be		
314	a big deal in the long-term & prompt engineer is not		
315	the job of the future AI gets easier. You can already		
316	see in Midjourney how basic prompts went from com-		
317	plex in v3 to easy in v4. Same with ChatGPT to Bing.		
318	https://t.co/BTtSN4oVF4 .		
319	Konstantin Hebenstreit, Robert Praas, Louis P.		
320	Kiesewetter, and Matthias Samwald. 2023.		
321	An automatically discovered chain-of-thought		
322	prompt generalizes to novel models and datasets .		
323	ArXiv:2305.02897 [cs] .		