Prompt Engineering Haw-Shiuan Chang

LLM Development

Internet low-quality text (e.g., from trolls or haters)



How do you improve LLM's performance without training?

- Architectures
 - MLP ullet
 - RNN lacksquare
 - Transformer \bullet
- Training Stages
 - Pretraining \bullet
 - Supervised Fine-tuning (SFT) \bullet
 - Alignment \bullet
 - Learning from Human Feedback (LHF)
 - Reasoning \bullet

Available Methods

Training Data

Reward Sources

Human Annotation

Positive + Negative

LLM as Judge

Only Positive-

None

None

None

Reward Usage



Prompt Engineering

Easy

Inference-time Improvement

- Prompt engineering
 - In-context learning
- Decoding
- Agentic
 - RAG
 - Tools
 - Assistant
 - Multi-LLM collaboration



Human brain is also almost a black box



https://blog.ml.cmu.edu/2019/05/17/explaining-a-black-boxusing-deep-variational-information-bottleneck-approach/





Prompt Engineering

- What prompts are better?
- Chain of Thought
- - Similar to Best of N
- In-context Learning

Self-Consistency and Tree of Thoughts / Beam Search

Which Prompts are better?



Demystifying Prompts in Language Models via Perplexity Estimation (<u>https://arxiv.org/pdf/2212.04037</u>)

Which Prompt is better?

- User: Please generate a sci-fi story
- LLM: OOXX
- User: Please revise the story to reveal a big secret of the main character
- LLM: XXOO
- User: Please revise the story to ...

- User: Please generate a sci-fi story
- Constraint 1: Please revise the story to reveal a big secret of the main character
- Constraint 2: Please revise the story to ...
- LLM: XXOO



Chain of Thoughts

Our experiments on CoT improvements Soft Reasoning Mathematica 0.8 0.2 Commonsense Symbolic Knowledge Zero-shot direct answer Zero-shot CoT



CoT Performance Improvement Across Tasks Aggregated by Paper and Category

Figure 2: Results from our meta-analysis (grey dots) aggregated by paper and category (blue dots).

Self Consistency and Tree of Thoughts



Figure 1: Schematic illustrating various approaches to problem solving with LLMs. Each rectangle box represents a *thought*, which is a coherent language sequence that serves as an intermediate step toward problem solving. See concrete examples of how thoughts are generated, evaluated, and searched in Figures 2,4,6.





The reward model could be anything. For example, LM probability (beam search), answer quality scorer, profanity/toxicity filter, sentiment classifier, PRM

Best of N

Best-of-N vs Beam Search



Figure 2 | Comparing different PRM search methods. Left: Best-of-N samples N full answers and then selects the best answer according to the PRM final score. Center: Beam search samples N candidates at each step, and selects the top M according to the PRM to continue the search from. Right: lookahead-search extends each step in beam-search to utilize a k-step lookahead while assessing which steps to retain and continue the search from. Thus lookahead-search needs more compute.

Scaling LLM Test-Time Compute Optimally can be More Effective than Scaling Model Parameters (<u>https://arxiv.org/pdf/2408.03314</u>)

Usefulness of Guidance Depends on the Difficulty

Comparing Beam Search and Best-of-N by Difficulty Level



• For LLM rather than LRM

In-context learning: by just conditioning on a few demonstrations of the task in its prefix

Circulation revenue has increased by 5% in Finland. // Positive

Panostaja did not disclose the purchase price. // Neutral

Paying off the national debt will be extremely painful. // Negative

The company anticipated its operating profit to improve. // _____



LLMs can solve novel downstream tasks

- Circulation revenue has increased by 5% in Finland. // Finance
- They defeated ... in the NFC Championship Game. // Sports
 - Apple ... development of in-house chips. // Tech
 - The company anticipated its operating profit to improve. // _____



https://ai.stanford.edu/blog/understanding-incontext/

Which Examples should we choose?

- Principle \bullet
 - Diversity \bullet
 - Quality \bullet
 - Relevancy/Coverage to query
- The optimal strategy seems to be task-dependent

In-context Learning with Retrieved Demonstrations for Language Models: A Survey (https://arxiv.org/pdf/2401.11624v1)

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



Representative Demonstration Selection for In-Context Learning with Two-Stage Determinantal Point Process (<u>https://</u> aclanthology.org/2023.emnlp-main.331.pdf)



What's in a demonstration?

Demonstrations

Circulation revenue has increased by 5% in Finland.

Panostaja did not disclose the purchase price.

Paying off the national debt will be extremely painful.

Test example

The acquisition will have an immediate positive impact. \n



How important is input-label mapping?



What about other aspects of the demonstrations?

Demos w/ gold labels	(Format ✓ Input distribution Circulation revenue has incr Panostaja did not disclose th
Demos w/ random labels	(Format ✓ Input distribution Circulation revenue has incr Panostaja did not disclose th
OOD Demos w/ random labels	(Format ✓ Input distribution Colour-printed lithograph. V Many accompanying market
Demos w/ random English words	(Format ✓ Input distribution Circulation revenue has incr Panostaja did not disclose th
Demos w/o labels	(Format X Input distribution Circulation revenue has incr Panostaja did not disclose th
Demos labels only	(Format X Input distribution positive neutral

n ✓ Label space ✓ Input-label mapping ✓) reased by 5% in Finland and 4% in Sweden in 2008. \n positive he purchase price. \n neutral

n ✓ Label space ✓ Input-label mapping ×) reased by 5% in Finland and 4% in Sweden in 2008. \n neutral ne purchase price. \n negative

n X Label space ✓ Input-label mapping X) Very good condition. Image size: 15 x 23 1/2 inches. \n neutral ting claims of cannabis products are often well-meaning. \n negative

n ✓ Label space × Input-label mapping ×) reased by 5% in Finland and 4% in Sweden in 2008. \n unanimity ne purchase price. \n wave

Label space X Input-label mapping X)
reased by 5% in Finland and 4% in Sweden in 2008.
ne purchase price.

X Label space ✓ Input-label mapping X)

What about other aspects of the demonstrations?





F: Format L: Label space I: Input distribution M: Input-Label Mapping

Followup work tells a different story:

- Yoo et al., 2022 shows that input-label different experimental conditions and eval metrics
- too much

A Survey to Recent Progress Towards Understanding In-Context Learning (<u>https://arxiv.org/pdf/2402.02212</u>)

mappings matter quite significantly when using

 Madaan and Yazdanbakhsh (2022) show that random rationales degrade chain-of-thought performance, but other modifications to the rationale (e.g., wrong equations) don't affect it

Question

- less important?
- Why?

- Prompt engineering becomes less important •
- For familiar prompts, the LLMs might output good answers
- For unfamiliar prompts, the LLMs might output good answers or bad answers
 - Bad answers will be filtered out in SFT/RLHF

After SFT and RLHF, does prompt engineering become more important or