

Prompt engineering

Prompt engineering is the process of structuring an instruction that can be interpreted and understood by a generative AI model. A prompt is natural language text describing the task that an AI should perform.

A prompt for a text-to-text <u>language model</u> can be a query such as "what is <u>Fermat's little theorem?</u>", [4] a command such as "write a poem about leaves falling", [5] or a longer statement including context, instructions, [6] and conversation history. Prompt engineering may involve phrasing a query, specifying a style, [5] providing relevant context [7] or assigning a role to the AI such as "Act as a native French speaker". [8] A prompt may include a few examples for a model to learn from, such as asking the model to complete "maison \rightarrow house, chat \rightarrow cat, chien \rightarrow " (the expected response being dog), [9] an approach called **few-shot learning**. [10]

When communicating with a <u>text-to-image</u> or a text-to-audio model, a typical prompt is a description of a desired output such as "a high-quality photo of an astronaut riding a horse" or "Lo-fi slow BPM electro chill with organic samples". Prompting a text-to-image model may involve adding, removing, emphasizing and re-ordering words to achieve a desired subject, style, 190 layout, lighting, and aesthetic.

In-context learning

Prompt engineering is enabled by **in-context learning**, defined as a model's ability to temporarily learn from prompts. The ability for in-context learning is an emergent ability 14 of large language models. Incontext learning itself is an emergent property of model scale, meaning breaks in downstream scaling laws occur such that its efficacy increases at a different rate in larger models than in smaller models.

In contrast to training and <u>fine-tuning</u> for each specific task, which are not temporary, what has been learnt during in-context learning is of a temporary nature. It does not carry the temporary contexts or biases, except the ones already present in the (pre)training <u>dataset</u>, from one conversation to the other. This result of "mesa-optimization" within <u>transformer</u> layers, is a form of <u>meta-learning</u> or "learning to learn". [21]

History

In 2018, researchers first proposed that all previously separate tasks in $\underline{\text{NLP}}$ could be cast as a question answering problem over a context. In addition, they trained a first single, joint, multi-task model that would answer any task-related question like "What is the sentiment" or "Translate this sentence to German" or "Who is the president?" [22]

In 2021, researchers fine-tuned one generatively pretrained model (T0) on performing 12 <u>NLP</u> tasks (using 62 datasets, as each task can have multiple datasets). The model showed good performance on new tasks, surpassing models trained directly on just performing one task (without pretraining). To solve a task, T0 is

given the task in a structured prompt, for example If $\{\{premise\}\}\$ is true, is it also true that $\{\{hypothesis\}\}$? ||| $\{\{entailed\}\}$. is the prompt used for making T0 solve entailment. [23]

A repository for prompts reported that over 2,000 public prompts for around 170 datasets were available in February 2022. [24]

In 2022 the *chain-of-thought* prompting technique was proposed by Google researchers. [17][25]

In 2023 several text-to-text and text-to-image prompt databases were publicly available. [26][27]

Text-to-text

Chain-of-thought

Chain-of-thought (CoT) prompting is a technique that allows <u>large language models</u> (LLMs) to solve a problem as a series of intermediate steps <u>before giving a final answer</u>. Chain-of-thought prompting improves reasoning ability by inducing the model to answer a multi-step problem with steps of reasoning that mimic a <u>train of thought</u>. It allows large language models to overcome difficulties with some reasoning tasks that require <u>logical thinking</u> and multiple steps to solve, such as <u>arithmetic</u> or <u>commonsense reasoning</u> questions. [31][32][33]

For example, given the question "Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?", a CoT prompt might induce the LLM to answer "A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9." [17]

As originally proposed, [17] each CoT prompt included a few Q&A examples. This made it a *few-shot* prompting technique. However, simply appending the words "Let's think step-by-step", [34] has also proven effective, which makes CoT a *zero-shot* prompting technique. This allows for better scaling as a user no longer needs to formulate many specific CoT Q&A examples. [35]

When applied to \underline{PaLM} , a 540B parameter $\underline{language\ model}$, CoT prompting significantly aided the model, allowing it to perform comparably with task-specific $\underline{fine\text{-tuned}}$ models on several tasks, achieving $\underline{state\ of\ mathematical\ reasoning\ benchmark}$. It is possible to fine-tune models on CoT reasoning datasets to enhance this capability further and stimulate better $\underline{fine\ models}$.

Example:[34]

```
Q: {question}
A: Let's think step by step.
```

Other techniques

Chain-of-thought prompting is just one of many prompt-engineering techniques. Various other techniques have been proposed. At least 29 distinct techniques have been published. [38]

Chain-of-Symbol (CoS) Prompting

Chain-of-Symbol prompting in conjunction with CoT prompting assists LLMs with its difficulty of spatial reasoning in text. In other words, using random symbols such as ' / ' assist the LLM to interpret spacing in text. This assists in reasoning and increases the performance of the LLM. [39]

Example:[39]

```
Input:

There are a set of bricks. The yellow brick C is on top of the brick E. The yellow brick D is on top of the brick A. The yellow brick E is on top of the brick D. The white brick A is on top of the brick B. For the brick B, the color is white. Now we have to get a specific brick. The bricks must now be grabbed from top to bottom, and if the lower brick is to be grabbed, the upper brick must be removed first. How to get brick D?

B/A/D/E/C
C/E
E/D
D
Output:

So we get the result as C, E, D.
```

Generated knowledge prompting

Generated knowledge $prompting^{[40]}$ first prompts the model to generate relevant facts for completing the prompt, then proceed to complete the prompt. The completion quality is usually higher, as the model can be conditioned on relevant facts.

Example: [40]

```
Generate some knowledge about the concepts in the input.
Input: {question}
Knowledge:
```

Least-to-most prompting

Least-to-most prompting [41] prompts a model to first list the sub-problems to a problem, then solve them in sequence, such that later sub-problems can be solved with the help of answers to previous sub-problems.

Example: [41]

```
Q: {question}
A: Let's break down this problem:
1.
```

Self-consistency decoding

Self-consistency decoding [42] performs several chain-of-thought rollouts, then selects the most commonly reached conclusion out of all the rollouts. If the rollouts disagree by a lot, a human can be queried for the correct chain of thought. [43]

Complexity-based prompting

Complexity-based prompting^[44] performs several CoT rollouts, then select the rollouts with the longest chains of thought, then select the most commonly reached conclusion out of those.

Self-refine

Self-refine^[45] prompts the LLM to solve the problem, then prompts the LLM to critique its solution, then prompts the LLM to solve the problem again in view of the problem, solution, and critique. This process is repeated until stopped, either by running out of tokens, time, or by the LLM outputting a "stop" token.

Example critique: [45]

```
I have some code. Give one suggestion to improve readability. Don't fix the code, just
give a suggestion.
Code: {code}
Suggestion:
```

Example refinement:

```
Code: {code}
Let's use this suggestion to improve the code.
Suggestion: {suggestion}
New Code:
```

Tree-of-thought

Tree-of-thought prompting [46] generalizes chain-of-thought by prompting the model to generate one or more "possible next steps", and then running the model on each of the possible next steps by <u>breadth-first</u>, <u>beam</u>, or some other method of tree search. [47]

Maieutic prompting

<u>Maieutic</u> prompting is similar to tree-of-thought. The model is prompted to answer a question with an explanation. The model is then prompted to explain parts of the explanation, and so on. Inconsistent explanation trees are pruned or discarded. This improves performance on complex commonsense reasoning. [48]

Example: [48]

```
Q: {question}
A: True, because

Q: {question}
A: False, because
```

Directional-stimulus prompting

Directional-stimulus prompting [49] includes a hint or cue, such as desired keywords, to guide a language model toward the desired output.

Example:^[49]

```
Article: {article}
Keywords:

Article: {article}
Q: Write a short summary of the article in 2-4 sentences that accurately incorporates the provided keywords.
Keywords: {keywords}
A:
```

Prompting to disclose uncertainty

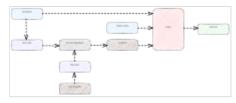
By default, the output of language models may not contain estimates of uncertainty. The model may output text that appears confident, though the underlying token predictions have low <u>likelihood</u> scores. Large language models like <u>GPT-4</u> can have accurately <u>calibrated</u> likelihood scores in their token predictions, [50] and so the model output uncertainty can be directly estimated by reading out the token prediction likelihood scores.

But if one cannot access such scores (such as when one is accessing the model through a restrictive API), uncertainty can still be estimated and incorporated into the model output. One simple method is to prompt the model to use words to estimate uncertainty. Another is to prompt the model to refuse to answer in a standardized way if the input does not satisfy conditions.

Automatic prompt generation

Retrieval-augmented generation

Retrieval-Augmented Generation (RAG) is a two-phase process involving document retrieval and answer formulation by a Large Language Model (LLM). The initial phase utilizes dense embeddings to retrieve documents. This retrieval can be based on a variety of database formats depending on the use case, such as a vector database, summary index, tree index, or keyword table index. [51]



Two-phase process of document retrieval using dense embeddings and Large Language Model (LLM) for answer formulation

In response to a query, a document retriever selects the most relevant documents. This relevance is typically determined by first

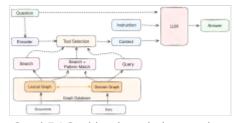
encoding both the query and the documents into vectors, then identifying documents whose vectors are closest in Euclidean distance to the query vector. Following document retrieval, the LLM generates an output that incorporates information from both the query and the retrieved documents. This method is particularly beneficial for handling proprietary or dynamic information that was not included in the initial training or fine-tuning phases of the model. RAG is also notable for its use of "few-shot" learning, where the model uses a small number of examples, often automatically retrieved from a database, to inform its outputs.

Graph retrieval-augmented generation

GraphRAG,^[53] coined by Microsoft Research, extends RAG such that instead of relying solely on vector similarity (as in most RAG approaches), GraphRAG uses the LLM-generated knowledge graph. This graph allows the model to connect disparate pieces of information, synthesize insights, and holistically

understand summarized semantic concepts over large data collections.

Researchers have demonstrated GraphRAG's effectiveness using datasets like the Violent Incident Information from News Articles (VIINA). By combining LLM-generated knowledge graphs with graph machine learning, GraphRAG substantially improves both the comprehensiveness and diversity of generated answers for global sensemaking questions.



GraphRAG with a knowledge graph combining access patterns for unstructured, structured and mixed data.

Earlier work showed the effectiveness of using a <u>knowledge graph</u> for question answering using text-to-query generation. [55] These

techniques can be combined to perform search across both unstructured and structured data, providing expanded context and improved ranking.

Using language models to generate prompts

Large language models (LLM) themselves can be used to compose prompts for large language models. [56][57][58]

The *automatic prompt engineer* algorithm uses one LLM to <u>beam search</u> over prompts for another LLM: [59]

- There are two LLMs. One is the target LLM, and another is the prompting LLM.
- Prompting LLM is presented with example input-output pairs, and asked to generate instructions that could have caused a model following the instructions to generate the outputs, given the inputs.
- Each of the generated instructions is used to prompt the target LLM, followed by each of the inputs. The log-probabilities of the outputs are computed and added. This is the score of the instruction.
- The highest-scored instructions are given to the prompting LLM for further variations.
- Repeat until some stopping criteria is reached, then output the highest-scored instructions.

CoT examples can be generated by LLM themselves. In "auto-CoT", [60] a library of questions are converted to vectors by a model such as <u>BERT</u>. The question vectors are <u>clustered</u>. Questions nearest to the centroids of each cluster are selected. An LLM does zero-shot CoT on each question. The resulting CoT examples are added to the dataset. When prompted with a new question, CoT examples to the nearest questions can be retrieved and added to the prompt.

Text-to-image

In 2022, <u>text-to-image</u> models like <u>DALL-E 2</u>, <u>Stable Diffusion</u>, and <u>Midjourney</u> were released to the public. These models take text prompts as input and use them to generate <u>AI art</u> images. Text-to-image models typically do not understand grammar and sentence structure in the same way as <u>large language</u> models, and require a different set of prompting techniques.

Prompt formats

A text-to-image prompt commonly includes a description of the subject of the art (such as *bright orange poppies*), the desired medium (such as *digital painting* or *photography*), style (such as *hyperrealistic* or *pop-art*), lighting (such as *rim lighting* or *crepuscular rays*), color and texture. [63]

The <u>Midjourney</u> documentation encourages short, descriptive prompts: instead of "Show me a picture of lots of blooming California poppies, make them bright, vibrant orange, and draw them in an illustrated style with colored pencils", an effective prompt might be "Bright orange California poppies drawn with colored pencils". [62]

Word order affects the output of a text-to-image prompt. Words closer to the start of a prompt may be emphasized more heavily. [1]

Artist styles

Some text-to-image models are capable of imitating the style of particular artists by name. For example, the phrase *in the style of Greg Rutkowski* has been used in Stable Diffusion and Midjourney prompts to generate images in the distinctive style of Polish digital artist Greg Rutkowski. [64]

Negative prompts



Demonstration of the effect of negative prompts on images generated with Stable Diffusion

Top: no negative prompt

Centre: "green trees"

■ **Bottom**: "round stones, round rocks"

Text-to-image models do not natively understand negation. The prompt "a party with no cake" is likely to produce an image including a cake. As an alternative, *negative prompts* allow a user to indicate, in a separate prompt, which terms should **not** appear in the resulting image. A common approach is to include generic undesired terms such as *ugly*, *boring*, *bad anatomy* in the negative prompt for an image.

Text-to-video

<u>Text-to-video</u> (TTV) generation is an emerging technology enabling the creation of videos directly from textual descriptions. This field holds potential for transforming video production, animation, and storytelling. By utilizing the power of artificial intelligence, TTV allows users to bypass traditional video editing tools and translate their ideas into moving images.

Models include:

- Runway Gen-2 Offers a user-friendly interface and supports various video styles
- Lumiere Designed for high-resolution video generation [66]
- Make-a-Video Focuses on creating detailed and diverse video outputs^[67]
- OpenAl's Sora As yet unreleased, Sora purportedly can produce high-resolution videos^{[68][69]}

Non-text prompts

Some approaches augment or replace natural language text prompts with non-text input.

Textual inversion and embeddings

For text-to-image models, "Textual inversion" performs an optimization process to create a new <u>word embedding</u> based on a set of example images. This embedding vector acts as a "pseudo-word" which can be included in a prompt to express the content or style of the examples.

Image prompting

In 2023, <u>Meta</u>'s AI research released Segment Anything, a <u>computer vision</u> model that can perform <u>image</u> <u>segmentation</u> by prompting. As an alternative to text prompts, Segment Anything can accept bounding boxes, segmentation masks, and foreground/background points. [71]

Using gradient descent to search for prompts

In "prefix-tuning", [72] "prompt tuning" or "soft prompting", [73] floating-point-valued vectors are searched directly by gradient descent, to maximize the log-likelihood on outputs.

Formally, let $\mathbf{E} = \{e_1, \dots, e_k\}$ be a set of soft prompt tokens (tunable embeddings), while $\mathbf{X} = \{\mathbf{x_1}, \dots, \mathbf{x_m}\}$ and $\mathbf{Y} = \{\mathbf{y_1}, \dots, \mathbf{y_n}\}$ be the token embeddings of the input and output respectively. During training, the tunable embeddings, input, and output tokens are concatenated into a single sequence $\mathbf{concat}(\mathbf{E}; \mathbf{X}; \mathbf{Y})$, and fed to the large language models (LLM). The losses are computed

over the \mathbf{Y} tokens; the gradients are <u>backpropagated</u> to prompt-specific parameters: in prefix-tuning, they are parameters associated with the prompt tokens at each layer; in prompt tuning, they are merely the soft tokens added to the vocabulary. [74]

More formally, this is prompt tuning. Let an LLM be written as LLM(X) = F(E(X)), where X is a sequence of linguistic tokens, E is the token-to-vector function, and F is the rest of the model. In prefixtuning, one provide a set of input-output pairs $\{(X^i,Y^i)\}_i$, and then use gradient descent to search for $\arg\max_{\tilde{Z}}\sum_i \log Pr[Y^i|\tilde{Z}*E(X^i)]$. In words, $\log Pr[Y^i|\tilde{Z}*E(X^i)]$ is the log-likelihood of outputting Y^i , if the model first encodes the input X^i into the vector $E(X^i)$, then prepend the vector with the "prefix vector" \tilde{Z} , then apply F.

For prefix tuning, it is similar, but the "prefix vector" $\tilde{\boldsymbol{Z}}$ is preappended to the hidden states in every layer of the model.

An earlier result [75] uses the same idea of gradient descent search, but is designed for masked language models like BERT, and searches only over token sequences, rather than numerical vectors. Formally, it searches for $\underset{\tilde{X}}{\operatorname{arg}} \max_{\tilde{X}} \sum_{i} \log Pr[Y^{i}|\tilde{X}*X^{i}]$ where \tilde{X} is ranges over token sequences of a specified length.

Prompt injection

Prompt injection is a family of related <u>computer security exploits</u> carried out by getting a <u>machine learning</u> model (such as an LLM) which was trained to follow human-given instructions to follow instructions provided by a malicious user. This stands in contrast to the intended operation of instruction-following systems, wherein the ML model is intended only to follow trusted instructions (prompts) provided by the ML model's operator. [76][77][78]

See also

Social engineering (security)

References

- 1. Diab, Mohamad; Herrera, Julian; Chernow, Bob (2022-10-28). "Stable Diffusion Prompt Book" (https://cdn.openart.ai/assets/Stable%20Diffusion%20Prompt%20Book%20From%20 OpenArt%2011-13.pdf) (PDF). Retrieved 2023-08-07. "Prompt engineering is the process of structuring words that can be interpreted and understood by a *text-to-image* model. Think of it as the language you need to speak in order to tell an AI model what to draw."
- 2. Ziegler, Albert; Berryman, John (17 July 2023). "A developer's guide to prompt engineering and LLMs" (https://github.blog/2023-07-17-prompt-engineering-guide-generative-ai-llms/). The GitHub Blog. "Prompt engineering is the art of communicating with a generative Al model."

- 3. Radford, Alec; Wu, Jeffrey; Child, Rewon; Luan, David; Amodei, Dario; Sutskever, Ilya (2019). "Language Models are Unsupervised Multitask Learners" (https://cdn.openai.com/bet ter-language-models/language_models_are_unsupervised_multitask_learners.pdf) (PDF). OpenAI. "We demonstrate language models can perform down-stream tasks in a zero-shot setting without any parameter or architecture modification"
- 4. "Introducing ChatGPT" (https://openai.com/blog/chatgpt). OpenAl Blog. 2022-11-30. Retrieved 2023-08-16. "what is the fermat's little theorem"
- 5. Robinson, Reid (August 3, 2023). "How to write an effective GPT-3 or GPT-4 prompt" (https://zapier.com/blog/gpt-prompt/). Zapier. Retrieved 2023-08-14. " "Basic prompt: 'Write a poem about leaves falling.' Better prompt: 'Write a poem in the style of Edgar Allan Poe about leaves falling.' "
- 6. Gouws-Stewart, Natasha (June 16, 2023). <u>"The ultimate guide to prompt engineering your GPT-3.5-Turbo model" (https://masterofcode.com/blog/the-ultimate-guide-to-gpt-prompt-engineering)</u>. *masterofcode.com*.
- 7. Greenberg, J., Laura (31 May 2023). "How to Prime and Prompt ChatGPT for More Reliable Contract Drafting Support" (https://contractnerds.com/how-to-prime-and-prompt-chatgpt-for-more-reliable-contract-drafting-support). contractnerds.com. Retrieved 24 July 2023.
- 8. "GPT Best Practices" (https://platform.openai.com/docs/guides/gpt-best-practices). OpenAl. Retrieved 2023-08-16.
- Garg, Shivam; Tsipras, Dimitris; Liang, Percy; Valiant, Gregory (2022). "What Can Transformers Learn In-Context? A Case Study of Simple Function Classes". arXiv:2208.01066 (https://arxiv.org/abs/2208.01066) [cs.CL (https://arxiv.org/archive/cs.CL)].
- 10. Brown, Tom; Mann, Benjamin; Ryder, Nick; Subbiah, Melanie; Kaplan, Jared D.; Dhariwal, Prafulla; Neelakantan, Arvind (2020). "Language models are few-shot learners". *Advances in Neural Information Processing Systems*. **33**: 1877–1901. arXiv:2005.14165 (https://arxiv.org/abs/2005.14165).
- 11. Heaven, Will Douglas (April 6, 2022). "This horse-riding astronaut is a milestone on Al's long road towards understanding" (https://www.technologyreview.com/2022/04/06/1049061/dalle-openai-gpt3-ai-agi-multimodal-image-generation/). MIT Technology Review. Retrieved 2023-08-14.
- 12. Wiggers, Kyle (2023-06-12). "Meta open sources an Al-powered music generator" (https://tec hcrunch.com/2023/06/12/meta-open-sources-an-ai-powered-music-generator/). TechCrunch. Retrieved 2023-08-15. "Next, I gave a more complicated prompt to attempt to throw MusicGen for a loop: "Lo-fi slow BPM electro chill with organic samples." "
- 13. "How to Write AI Photoshoot Prompts: A Guide for Better Product Photos" (https://claid.ai/blog/article/prompt-guide/). claid.ai. June 12, 2023. Retrieved June 12, 2023.
- 14. Wei, Jason; Tay, Yi; Bommasani, Rishi; Raffel, Colin; Zoph, Barret; Borgeaud, Sebastian; Yogatama, Dani; Bosma, Maarten; Zhou, Denny; Metzler, Donald; Chi, Ed H.; Hashimoto, Tatsunori; Vinyals, Oriol; Liang, Percy; Dean, Jeff; Fedus, William (31 August 2022). "Emergent Abilities of Large Language Models". arXiv:2206.07682 (https://arxiv.org/abs/2206.07682) [cs.CL (https://arxiv.org/archive/cs.CL)]. "In prompting, a pre-trained language model is given a prompt (e.g. a natural language instruction) of a task and completes the response without any further training or gradient updates to its parameters... The ability to perform a task via few-shot prompting is emergent when a model has random performance until a certain scale, after which performance increases to well-above random"
- 15. Caballero, Ethan; Gupta, Kshitij; Rish, Irina; Krueger, David (2022). "Broken Neural Scaling Laws". International Conference on Learning Representations (ICLR), 2023.

- 16. Wei, Jason; Tay, Yi; Bommasani, Rishi; Raffel, Colin; Zoph, Barret; Borgeaud, Sebastian; Yogatama, Dani; Bosma, Maarten; Zhou, Denny; Metzler, Donald; Chi, Ed H.; Hashimoto, Tatsunori; Vinyals, Oriol; Liang, Percy; Dean, Jeff; Fedus, William (31 August 2022). "Emergent Abilities of Large Language Models". arXiv:2206.07682 (https://arxiv.org/abs/220 6.07682) [cs.CL (https://arxiv.org/archive/cs.CL)].
- 17. Wei, Jason; Wang, Xuezhi; Schuurmans, Dale; Bosma, Maarten; Ichter, Brian; Xia, Fei; Chi, Ed H.; Le, Quoc V.; Zhou, Denny (31 October 2022). *Chain-of-Thought Prompting Elicits Reasoning in Large Language Models* (https://proceedings.neurips.cc/paper_files/paper/202 2/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html). Advances in Neural Information Processing Systems (NeurIPS 2022). Vol. 35. arXiv:2201.11903 (https://arxiv.org/abs/2201.11903).
- 18. Musser, George. "How Al Knows Things No One Told It" (https://www.scientificamerican.co m/article/how-ai-knows-things-no-one-told-it/). <u>Scientific American</u>. Retrieved 17 May 2023. "By the time you type a query into ChatGPT, the network should be fixed; unlike humans, it should not continue to learn. So it came as a surprise that LLMs do, in fact, learn from their users' prompts—an ability known as in-context learning."
- 19. Johannes von Oswald; Niklasson, Eyvind; Randazzo, Ettore; Sacramento, João; Mordvintsev, Alexander; Zhmoginov, Andrey; Vladymyrov, Max (2022). "Transformers learn in-context by gradient descent". arXiv:2212.07677 (https://arxiv.org/archive/cs.LG)]. "Thus we show how trained Transformers become mesaoptimizers i.e. learn models by gradient descent in their forward pass"
- 20. "Mesa-Optimization" (https://www.alignmentforum.org/tag/mesa-optimization). 31 May 2019. Retrieved 17 May 2023. "Mesa-Optimization is the situation that occurs when a learned model (such as a neural network) is itself an optimizer."
- 21. Garg, Shivam; Tsipras, Dimitris; Liang, Percy; Valiant, Gregory (2022). "What Can Transformers Learn In-Context? A Case Study of Simple Function Classes". arXiv:2208.01066 (https://arxiv.org/abs/2208.01066) [cs.CL (https://arxiv.org/archive/cs.CL)]. "Training a model to perform in-context learning can be viewed as an instance of the more general learning-to-learn or meta-learning paradigm"
- 22. McCann, Bryan; Shirish, Nitish; Xiong, Caiming; Socher, Richard (2018). "The Natural Language Decathlon: Multitask Learning as Question Answering". arXiv:1806.08730 (https://arxiv.org/archive/cs.CL)].
- 23. Sanh, Victor; et al. (2021). "Multitask Prompted Training Enables Zero-Shot Task Generalization". arXiv:2110.08207 (https://arxiv.org/abs/2110.08207) [cs.LG (https://arxiv.org/archive/cs.LG)].
- 24. Bach, Stephen H.; Sanh, Victor; Yong, Zheng-Xin; Webson, Albert; Raffel, Colin; Nayak, Nihal V.; Sharma, Abheesht; Kim, Taewoon; M Saiful Bari; Fevry, Thibault; Alyafeai, Zaid; Dey, Manan; Santilli, Andrea; Sun, Zhiqing; Ben-David, Srulik; Xu, Canwen; Chhablani, Gunjan; Wang, Han; Jason Alan Fries; Al-shaibani, Maged S.; Sharma, Shanya; Thakker, Urmish; Almubarak, Khalid; Tang, Xiangru; Radev, Dragomir; Mike Tian-Jian Jiang; Rush, Alexander M. (2022). "PromptSource: An Integrated Development Environment and Repository for Natural Language Prompts". arXiv:2202.01279 (https://arxiv.org/abs/2202.01279) [cs.LG (https://arxiv.org/archive/cs.LG)].
- 25. Wei, Jason; Zhou (11 May 2022). "Language Models Perform Reasoning via Chain of Thought" (https://ai.googleblog.com/2022/05/language-models-perform-reasoning-via.html). ai.googleblog.com. Retrieved 10 March 2023.
- 26. Chen, Brian X. (2023-06-23). "How to Turn Your Chatbot Into a Life Coach" (https://www.nytimes.com/2023/06/23/technology/ai-chatbot-life-coach.html). The New York Times.

- 27. Chen, Brian X. (2023-05-25). "Get the Best From ChatGPT With These Golden Prompts" (htt ps://www.nytimes.com/2023/05/25/technology/ai-chatbot-chatgpt-prompts.html). *The New York Times*. ISSN 0362-4331 (https://www.worldcat.org/issn/0362-4331). Retrieved 2023-08-16.
- 28. McAuliffe, Zachary. "Google's Latest Al Model Can Be Taught How to Solve Problems" (http s://www.cnet.com/tech/services-and-software/googles-latest-ai-model-can-be-taught-how-to-solve-problems/). CNET. Retrieved 10 March 2023. " 'Chain-of-thought prompting allows us to describe multistep problems as a series of intermediate steps,' Google CEO Sundar Pichai"
- 29. McAuliffe, Zachary. "Google's Latest Al Model Can Be Taught How to Solve Problems" (http s://www.cnet.com/tech/services-and-software/googles-latest-ai-model-can-be-taught-how-to-solve-problems/). CNET. Retrieved 10 March 2023.
- 30. Sharan Narang and Aakanksha Chowdhery (2022-04-04). "Pathways Language Model (PaLM): Scaling to 540 Billion Parameters for Breakthrough Performance" (https://ai.googleblog.com/2022/04/pathways-language-model-palm-scaling-to.html).
- 31. Dang, Ekta (8 February 2023). <u>"Harnessing the power of GPT-3 in scientific research" (http s://venturebeat.com/ai/harnessing-the-power-of-gpt-3-in-scientific-research/)</u>. *VentureBeat*. Retrieved 10 March 2023.
- 32. Montti, Roger (13 May 2022). "Google's Chain of Thought Prompting Can Boost Today's Best Algorithms" (https://www.searchenginejournal.com/google-chain-of-thought-prompting/4 50106/). Search Engine Journal. Retrieved 10 March 2023.
- 33. Ray, Tiernan. "Amazon's Alexa scientists demonstrate bigger Al isn't always better" (https://www.zdnet.com/article/amazons-alexa-scientists-demonstrate-bigger-ai-isnt-always-better/).

 ZDNET. Retrieved 10 March 2023.
- 34. Kojima, Takeshi; Shixiang Shane Gu; Reid, Machel; Matsuo, Yutaka; Iwasawa, Yusuke (2022). "Large Language Models are Zero-Shot Reasoners". arXiv:2205.11916 (https://arxiv.org/archive/cs.CL)].
- 35. Dickson, Ben (30 August 2022). <u>"LLMs have not learned our language we're trying to learn theirs"</u> (https://venturebeat.com/ai/llms-have-not-learned-our-language-were-trying-to-learn-theirs%EF%BF%BC/). *VentureBeat*. Retrieved 10 March 2023.
- 36. Chung, Hyung Won; Hou, Le; Longpre, Shayne; Zoph, Barret; Tay, Yi; Fedus, William; Li, Yunxuan; Wang, Xuezhi; Dehghani, Mostafa; Brahma, Siddhartha; Webson, Albert; Gu, Shixiang Shane; Dai, Zhuyun; Suzgun, Mirac; Chen, Xinyun; Chowdhery, Aakanksha; Castro-Ros, Alex; Pellat, Marie; Robinson, Kevin; Valter, Dasha; Narang, Sharan; Mishra, Gaurav; Yu, Adams; Zhao, Vincent; Huang, Yanping; Dai, Andrew; Yu, Hongkun; Petrov, Slav; Chi, Ed H.; Dean, Jeff; Devlin, Jacob; Roberts, Adam; Zhou, Denny; Le, Quoc V.; Wei, Jason (2022). "Scaling Instruction-Finetuned Language Models". arXiv:2210.11416 (https://arxiv.org/abs/2210.11416) [cs.LG (https://arxiv.org/archive/cs.LG)].
- 37. Wei, Jason; Tay, Yi (29 November 2022). "Better Language Models Without Massive Compute" (https://ai.googleblog.com/2022/11/better-language-models-without-massive.htm l). ai.googleblog.com. Retrieved 10 March 2023.
- 38. Sahoo, Pranab; Singh, Ayush Kumar; Saha, Sriparna; Jain, Vinija; Mondal, Samrat; Chadha, Aman (2024-02-05), *A Systematic Survey of Prompt Engineering in Large Language Models: Techniques and Applications*, arXiv:2402.07927 (https://arxiv.org/abs/2402.07927)
- 39. Hu, Hanxu; Lu, Hongyuan; Zhang, Huajian; Song, Yun-Ze; Lam, Wai; Zhang, Yue (2023-10-03), Chain-of-Symbol Prompting Elicits Planning in Large Language Models, arXiv:2305.10276 (https://arxiv.org/abs/2305.10276)

- 40. Liu, Jiacheng; Liu, Alisa; Lu, Ximing; Welleck, Sean; West, Peter; Le Bras, Ronan; Choi, Yejin; Hajishirzi, Hannaneh (May 2022). "Generated Knowledge Prompting for Commonsense Reasoning" (https://aclanthology.org/2022.acl-long.225). Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Dublin, Ireland: Association for Computational Linguistics: 3154–3169. arXiv:2110.08387 (https://arxiv.org/abs/2110.08387). doi:10.18653/v1/2022.acl-long.225 (https://doi.org/10.18653%2Fv1%2F2022.acl-long.225). S2CID 239016123 (https://api.semanticscholar.org/CorpusID:239016123).
- 41. Zhou, Denny; Schärli, Nathanael; Hou, Le; Wei, Jason; Scales, Nathan; Wang, Xuezhi; Schuurmans, Dale; Cui, Claire; Bousquet, Olivier; Le, Quoc; Chi, Ed (2022-05-01). "Least-to-Most Prompting Enables Complex Reasoning in Large Language Models". arXiv:2205.10625 (https://arxiv.org/abs/2205.10625) [cs.Al (https://arxiv.org/archive/cs.Al)]. "...least-to-most prompting. The key idea in this strategy is to break down a complex problem into a series of simpler subproblems and then solve them in sequence."
- 42. Wang, Xuezhi; Wei, Jason; Schuurmans, Dale; Le, Quoc; Chi, Ed; Narang, Sharan; Chowdhery, Aakanksha; Zhou, Denny (2022-03-01). "Self-Consistency Improves Chain of Thought Reasoning in Language Models". arXiv:2203.11171 (https://arxiv.org/abs/2203.11171) [cs.CL (https://arxiv.org/archive/cs.CL)].
- 43. Diao, Shizhe; Wang, Pengcheng; Lin, Yong; Zhang, Tong (2023-02-01). "Active Prompting with Chain-of-Thought for Large Language Models". arXiv:2302.12246 (https://arxiv.org/archive/cs.CL)].
- 44. Fu, Yao; Peng, Hao; Sabharwal, Ashish; Clark, Peter; Khot, Tushar (2022-10-01). "Complexity-Based Prompting for Multi-Step Reasoning". arXiv:2210.00720 (https://arxiv.org/archive/cs.CL)].
- 45. Madaan, Aman; Tandon, Niket; Gupta, Prakhar; Hallinan, Skyler; Gao, Luyu; Wiegreffe, Sarah; Alon, Uri; Dziri, Nouha; Prabhumoye, Shrimai; Yang, Yiming; Gupta, Shashank; Prasad Majumder, Bodhisattwa; Hermann, Katherine; Welleck, Sean; Yazdanbakhsh, Amir (2023-03-01). "Self-Refine: Iterative Refinement with Self-Feedback". arXiv:2303.17651 (https://arxiv.org/abs/2303.17651) [cs.CL (https://arxiv.org/archive/cs.CL)].
- 46. Long, Jieyi (2023-05-15). "Large Language Model Guided Tree-of-Thought". arXiv:2305.08291 (https://arxiv.org/abs/2305.08291) [cs.AI (https://arxiv.org/archive/cs.AI)].
- 47. Yao, Shunyu; Yu, Dian; Zhao, Jeffrey; Shafran, Izhak; Griffiths, Thomas L.; Cao, Yuan; Narasimhan, Karthik (2023-05-17). "Tree of Thoughts: Deliberate Problem Solving with Large Language Models". arXiv:2305.10601 (https://arxiv.org/abs/2305.10601) [cs.CL (https://arxiv.org/archive/cs.CL)].
- 48. Jung, Jaehun; Qin, Lianhui; Welleck, Sean; Brahman, Faeze; Bhagavatula, Chandra; Le Bras, Ronan; Choi, Yejin (2022). "Maieutic Prompting: Logically Consistent Reasoning with Recursive Explanations". arXiv:2205.11822 (https://arxiv.org/abs/2205.11822) [cs.CL (https://arxiv.org/archive/cs.CL)].
- 49. Li, Zekun; Peng, Baolin; He, Pengcheng; Galley, Michel; Gao, Jianfeng; Yan, Xifeng (2023). "Guiding Large Language Models via Directional Stimulus Prompting". arXiv:2302.11520 (ht tps://arxiv.org/abs/2302.11520) [cs.CL (https://arxiv.org/archive/cs.CL)]. "The directional stimulus serves as hints or cues for each input query to guide LLMs toward the desired output, such as keywords that the desired summary should include for summarization."
- 50. OpenAI (2023-03-27). "GPT-4 Technical Report". arXiv:2303.08774 (https://arxiv.org/abs/230 3.08774) [cs.CL (https://arxiv.org/archive/cs.CL)]. [See Figure 8.]
- 51. "How Each Index Works LlamaIndex > v0.10.17" (https://docs.llamaindex.ai/en/v0.10.17/module_guides/indexing/index_guide.html). docs.llamaindex.ai. Retrieved 2024-04-08.

- 52. Lewis, Patrick; Perez, Ethan; Piktus, Aleksandra; Petroni, Fabio; Karpukhin, Vladimir; Goyal, Naman; Küttler, Heinrich; Lewis, Mike; Yih, Wen-tau; Rocktäschel, Tim; Riedel, Sebastian; Kiela, Douwe (2020). "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks" (https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df748 1e5-Abstract.html). Advances in Neural Information Processing Systems. 33. Curran Associates, Inc.: 9459–9474. arXiv:2005.11401 (https://arxiv.org/abs/2005.11401).
- 53. *GraphRAG: Unlocking LLM discovery on narrative private data* (https://www.microsoft.com/e n-us/research/blog/graphrag-unlocking-llm-discovery-on-narrative-private-data/), 2024
- 54. Edge, Darren; Trinh, Ha; Cheng, Newman; Bradley, Joshua; Chao, Alex; Mody, Apurva; Truitt, Steven; Larson, Jonathan (2024), From Local to Global: A Graph RAG Approach to Query-Focused Summarization, arXiv:2404.16130 (https://arxiv.org/abs/2404.16130)
- 55. A Benchmark to Understand the Role of Knowledge Graphs on Large Language Model's Accuracy for Question Answering on Enterprise SQL Databases, 2023, arXiv:2311.07509 (https://arxiv.org/abs/2311.07509)
- 56. Fernando, Chrisantha; Banarse, Dylan; Michalewski, Henryk; Osindero, Simon; Rocktäschel, Tim (2023). "Promptbreeder: Self-Referential Self-Improvement Via Prompt Evolution". arXiv:2309.16797 (https://arxiv.org/abs/2309.16797). {{cite journal}}: Cite journal requires | journal (help)
- 57. Pryzant, Reid; Iter, Dan; Li, Jerry; Lee, Yin Tat; Zhu, Chenguang; Zeng, Michael (2023). "Automatic Prompt Optimization with "Gradient Descent" and Beam Search". arXiv:2305.03495 (https://arxiv.org/abs/2305.03495). {{cite journal}}: Cite journal requires | journal = (help)
- 58. Guo, Qingyan; Wang, Rui; Guo, Junliang; Li, Bei; Song, Kaitao; Tan, Xu; Liu, Guoqing; Bian, Jiang; Yang, Yujiu (2023). "Connecting Large Language Models with Evolutionary Algorithms Yields Powerful Prompt Optimizers". arXiv:2309.08532 (https://arxiv.org/abs/230 9.08532). {{cite journal}}: Cite journal requires | journal = (help)
- 59. Zhou, Yongchao; Ioan Muresanu, Andrei; Han, Ziwen; Paster, Keiran; Pitis, Silviu; Chan, Harris; Ba, Jimmy (2022-11-01). "Large Language Models Are Human-Level Prompt Engineers". arXiv:2211.01910 (https://arxiv.org/abs/2211.01910) [cs.LG (https://arxiv.org/archive/cs.LG)].
- 60. Zhang, Zhuosheng; Zhang, Aston; Li, Mu; Smola, Alex (2022-10-01). "Automatic Chain of Thought Prompting in Large Language Models". arXiv:2210.03493 (https://arxiv.org/abs/221 0.03493) [cs.CL (https://arxiv.org/archive/cs.CL)].
- 61. Monge, Jim Clyde (2022-08-25). "Dall-E2 VS Stable Diffusion: Same Prompt, Different Results" (https://medium.com/mlearning-ai/dall-e2-vs-stable-diffusion-same-prompt-different-results-e795c84adc56). *MLearning.ai*. Retrieved 2022-08-31.
- 62. "Prompts" (https://docs.midjourney.com/docs/prompts). Retrieved 2023-08-14.
- 63. "Stable Diffusion prompt: a definitive guide" (https://stable-diffusion-art.com/prompt-guide/). 2023-05-14. Retrieved 2023-08-14.
- 64. Heikkilä, Melissa (2022-09-16). "This Artist Is Dominating Al-Generated Art and He's Not Happy About It" (https://www.technologyreview.com/2022/09/16/1059598/this-artist-is-dominating-ai-generated-art-and-hes-not-happy-about-it/). MIT Technology Review. Retrieved 2023-08-14.
- 65. Max Woolf (2022-11-28). <u>"Stable Diffusion 2.0 and the Importance of Negative Prompts for Good Results" (https://minimaxir.com/2022/11/stable-diffusion-negative-prompt/)</u>. Retrieved 2023-08-14.
- 66. "Lumiere Google Research" (https://lumiere-video.github.io/). Lumiere Google Research. Retrieved 2024-02-25.
- 67. "Introducing Make-A-Video: An AI system that generates videos from text" (https://ai.meta.com/blog/generative-ai-text-to-video/). ai.meta.com. Retrieved 2024-02-25.

- 68. "Video generation models as world simulators" (https://openai.com/research/video-generation-models-as-world-simulators). *openai.com*. Retrieved 2024-02-25.
- 69. Team, PromptSora. "Understanding OpenAl's Sora: A Revolutionary Leap | PromptSora: Discover Prompts and Videos for Sora from Open Al" (https://promptsora.com/blog/understan ding-openai-sora-a-revolutionary-leap). *PromptSora*. Retrieved 2024-02-25.
- 70. Gal, Rinon; Alaluf, Yuval; Atzmon, Yuval; Patashnik, Or; Bermano, Amit H.; Chechik, Gal; Cohen-Or, Daniel (2022). "An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion". arxiv:2208.01618 (https://arxiv.org/abs/2208.01618) [cs.CV (https://arxiv.org/archive/cs.CV)]. "Using only 3-5 images of a user-provided concept, like an object or a style, we learn to represent it through new "words" in the embedding space of a frozen text-to-image model."
- 71. Kirillov, Alexander; Mintun, Eric; Ravi, Nikhila; Mao, Hanzi; Rolland, Chloe; Gustafson, Laura; Xiao, Tete; Whitehead, Spencer; Berg, Alexander C.; Lo, Wan-Yen; Dollár, Piotr; Girshick, Ross (2023-04-01). "Segment Anything". arXiv:2304.02643 (https://arxiv.org/abs/23 04.02643) [cs.CV (https://arxiv.org/archive/cs.CV)].
- 72. Li, Xiang Lisa; Liang, Percy (2021). "Prefix-Tuning: Optimizing Continuous Prompts for Generation". *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. pp. 4582–4597. doi:10.18653/V1/2021.ACL-LONG.353 (https://doi.org/10.18653%2FV1%2F2021.ACL-LONG.353). S2CID 230433941 (https://api.semanticscholar.org/CorpusID:230433941). "In this paper, we propose prefix-tuning, a lightweight alternative to fine-tuning... Prefix-tuning draws inspiration from prompting"
- 73. Lester, Brian; Al-Rfou, Rami; Constant, Noah (2021). "The Power of Scale for Parameter-Efficient Prompt Tuning". *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. pp. 3045–3059. arXiv:2104.08691 (https://arxiv.org/abs/2104.08691). doi:10.18653/V1/2021.EMNLP-MAIN.243 (https://doi.org/10.18653%2FV1%2F2021. EMNLP-MAIN.243). S2CID 233296808 (https://api.semanticscholar.org/CorpusID:233296808). "In this work, we explore "prompt tuning," a simple yet effective mechanism for learning "soft prompts"...Unlike the discrete text prompts used by GPT-3, soft prompts are learned through back-propagation"
- 74. Sun, Simeng; Liu, Yang; Iter, Dan; Zhu, Chenguang; Iyyer, Mohit (2023). "How Does In-Context Learning Help Prompt Tuning?". arXiv:2302.11521 (https://arxiv.org/abs/2302.11521) [cs.CL (https://arxiv.org/archive/cs.CL)].
- 75. Shin, Taylor; Razeghi, Yasaman; Logan IV, Robert L.; Wallace, Eric; Singh, Sameer (November 2020). "AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts" (https://aclanthology.org/2020.emnlp-main.346). Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). Online: Association for Computational Linguistics. pp. 4222–4235. doi:10.18653/v1/2020.emnlp-main.346 (https://doi.org/10.18653%2Fv1%2F2020.emnlp-main.346). S2CID 226222232 (https://api.semanticscholar.org/CorpusID:226222232).
- 76. Willison, Simon (12 September 2022). "Prompt injection attacks against GPT-3" (http://simon willison.net/2022/Sep/12/prompt-injection/). simonwillison.net. Retrieved 2023-02-09.
- 77. Papp, Donald (2022-09-17). "What's Old Is New Again: GPT-3 Prompt Injection Attack Affects AI" (https://hackaday.com/2022/09/16/whats-old-is-new-again-gpt-3-prompt-injection-attack-affects-ai/). Hackaday. Retrieved 2023-02-09.
- 78. Vigliarolo, Brandon (19 September 2022). "GPT-3 'prompt injection' attack causes bot bad manners" (https://www.theregister.com/2022/09/19/in_brief_security/). www.theregister.com. Retrieved 2023-02-09.

