# Reasoning with Language Model Prompting: A Survey

**Shuofei Qiao**[1*], **Yixin Ou**[1*], **Ningyu Zhang**[1†], **Xiang Chen**[1], **Yunzhi Yao**[1],
**Shumin Deng**[4], **Chuanqi Tan**[3], **Fei Huang**[3], **Huajun Chen**[1,2†]

[1] Zhejiang University & AZFT Joint Lab for Knowledge Engine
[2] Donghai Laboratory [3] Alibaba Group [4] National University of Singapore

{shuofei,ouyixin,zhangningyu,xiang_chen,yyztodd,huajunsir}@zju.edu.cn
shumin@nus.edu.sg {chuanqi.tcq,f.huang}@alibaba-inc.com

## Abstract

Reasoning, as an essential ability for complex problem-solving, can provide back-end support for various real-world applications, such as medical diagnosis, negotiation, etc. This paper provides a comprehensive survey of cutting-edge research on reasoning with language model prompting. We introduce research works with comparisons and summaries and provide systematic resources to help beginners. We also discuss the potential reasons for emerging such reasoning abilities and highlight future research directions[1].

## 1 Introduction

Reasoning ability lies at the heart of human intelligence, yet in natural language processing (NLP), modern neural networks can hardly reason from what they are told or have already known (Duan et al., 2020; Wang et al., 2021; Bhargava and Ng, 2022). Fortunately, with the revolutionary development of pre-training (Brown et al., 2020; Chen et al., 2021; Chowdhery et al., 2022), scaling up the size of language models (LMs) has shown to confer a range of reasoning abilities, such as arithmetic (Wang et al., 2022e; Lewkowycz et al., 2022), commonsense (Jung et al., 2022; Liu et al., 2022b), symbolic (Zhou et al., 2023; Khot et al., 2023) reasoning. As shown in Figure 1, such abilities may be unlocked by prompting strategies (Liu et al., 2022d) (e.g., *chain-of-thought (CoT) prompting* (Wei et al., 2022b), *generated knowledge prompting* (Liu et al., 2022c)), which can dramatically narrow the gap between human and machine intelligence. Likewise, a vast amount of work has been proposed in the NLP community; however, these approaches, scattered among various tasks, have not been systematically reviewed and analyzed.

---

\* Equal Contribution.
† Corresponding Author.
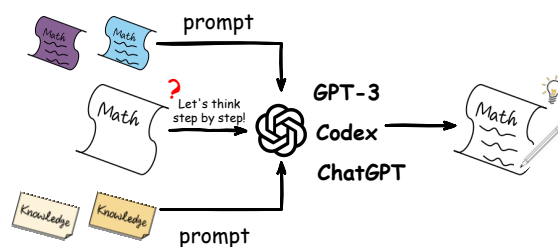[1]Resources are available at https://github.com/zjunlp/Prompt4ReasoningPapers (updated periodically).



Figure 1: Reasoning with language model prompting. In-context exemplars (colored ●, ●), knowledge (colored ●, ●) or just *Let's think step by step!* are as prompt to enhance language models reasoning.

**Organization of This Survey:** In this paper, we conduct the first survey of recent progress in reasoning with language model prompting. We first give some preliminaries on this direction (§2) and then propose to organize relevant works by taxonomy (§3). We further provide in-depth comparisons with discussion for insights (§4). To facilitate beginners who are interested in this field, we highlight some open resources (§5) as well as potential future directions (§6).

## 2 Preliminaries

In this section, we introduce preliminaries of reasoning with LM prompting. For standard prompting, given the reasoning question $\mathcal{Q}$, prompt $\mathcal{T}$ and parameterized probabilistic model $p_{\text{LM}}$, we aim to maximize the likelihood of answer $\mathcal{A}$ as:

$$p(\mathcal{A} \mid \mathcal{T}, \mathcal{Q}) = \prod_{i=1}^{|\mathcal{A}|} p_{\text{LM}}\left(a_i \mid \mathcal{T}, \mathcal{Q}, a_{<i}\right) \quad (1)$$

where $a_i$ and $|\mathcal{A}|$ denotes the $i$-th token and the length of the final answer respectively. For few-shot prompting, $\mathcal{T}$ is comprised of $\mathcal{K}$ exemplars of $(\mathcal{Q}, \mathcal{A})$ pair. CoT approaches further *add reasoning steps* $\mathcal{C}$ into prompt where $\mathcal{T} = \{(\mathcal{Q}_i, \mathcal{C}_i, \mathcal{A}_i)\}_{i=1}^{\mathcal{K}}$, thus Equation 1 can be reformed to:

$$p(\mathcal{A} \mid \mathcal{T}, \mathcal{Q}) = p\left(\mathcal{A} \mid \mathcal{T}, \mathcal{Q}, \mathcal{C}\right) p\left(\mathcal{C} \mid \mathcal{T}, \mathcal{Q}\right) \quad (2)$$

5368

Figure 2: Taxonomy tree

**Reasoning with Language Model Prompting**

- **Taxonomy of Methods (§3)**
  - **Strategy Enhanced Reasoning (§3.1)**
    - **Prompt Engineering (§3.1.1)**
      - **Single-Stage**: Contrastive (Paranjape et al., 2021), POTTER (Rajagopal et al., 2021), CoT (Wei et al., 2022b), ZeroCoT (Kojima et al., 2022), Complexity (Fu et al., 2023b), Multilingual (Shi et al., 2022), Auto-CoT (Zhang et al., 2023b), Table (Chen, 2022), AlgoPrompt (Zhou et al., 2022a), Active-Prompt (Diao et al., 2023), Automate-CoT (Shum et al., 2023)
      - **Multi-Stage**: iCAP (Wang et al., 2022a), SI (Creswell et al., 2022), Least-to-Most (Zhou et al., 2023), MAIEUTIC (Jung et al., 2022), Faithful (Creswell and Shanahan, 2022), Decomposed (Khot et al., 2023), Self-Ask (Press et al., 2022), Successive (Dua et al., 2022), LMLP (Zhang et al., 2022), LAMBADA (Kazemi et al., 2022), Iter-Decomp (Reppert et al., 2023)
    - **Process Optimization (§3.1.2)**
      - **Self-Optimization**: Calibrator (Ye and Durrett, 2022), Human-AI (Wiegreffe et al., 2022)
      - **Ensemble-Optimization**: Self-C (Wang et al., 2022e), DIVERSE (Li et al., 2022d), Complexity (Fu et al., 2023b), Self-V (Weng et al., 2022), MCR (Yoran et al., 2023)
      - **Iterative-Optimization**: STaR (Zelikman et al., 2022), LMSI (Huang et al., 2022), Reflexion (Shinn et al., 2023), Self-Refine (Madaan et al., 2023), REFINER (Paul et al., 2023)
    - **External Engine (§3.1.3)**
      - **Physical Simulator**: Mind's Eye (Liu et al., 2023)
      - **Code Interpreter**: COCOGEN (Madaan et al., 2022), PAL (Gao et al., 2022), PoT (Chen et al., 2022b), Faithful-CoT (Lyu et al., 2023), Versa-Decomp (Ye et al., 2023), SynPrompt (Shao et al., 2023), MathPrompter (Imani et al., 2023)
      - **Tool Learning**: Toolformer (Schick et al., 2023), ART (Paranjape et al., 2023), Chameleon (Lu et al., 2023a)
  - **Knowledge Enhanced Reasoning (§3.2)**
    - **Implicit Knowledge (§3.2.1)**: GenKnow (Liu et al., 2022c), RAINIER (Liu et al., 2022b), MT-CoT (Li et al., 2022b), PINTO (Wang et al., 2023), TSGP (Sun et al., 2022), DecompDistill (Shridhar et al., 2022), Teaching (Magister et al., 2022), Fine-tune-CoT (Ho et al., 2022), Specializing (Fu et al., 2023a)
    - **Explicit Knowledge (§3.2.2)**: LogicSolver (Yang et al., 2022b), Vote-$k$ (SU et al., 2023), PROMPTPG (Lu et al., 2023b), IRCoT (Trivedi et al., 2022), RR (He et al., 2023)
- **Taxonomy of Tasks (§5)**
  - **Arithmetic**: CoT (Wei et al., 2022b), Self-C (Wang et al., 2022e), Least-to-Most (Zhou et al., 2023), ZeroCoT (Kojima et al., 2022), Auto-CoT (Zhang et al., 2023b), LMSI (Huang et al., 2022), PAL (Gao et al., 2022), PoT (Chen et al., 2022b), Fine-tune-CoT (Ho et al., 2022)
  - **Commonsense**: CoT (Wei et al., 2022b), GenKnow (Liu et al., 2022c), Self-C (Wang et al., 2022e), Calibrator (Ye and Durrett, 2022), ZeroCoT (Kojima et al., 2022), Auto-CoT (Zhang et al., 2023b), COCOGEN (Madaan et al., 2022), LMSI (Huang et al., 2022), PINTO (Wang et al., 2023), RR (He et al., 2023)
  - **Logical**: Faithful (Creswell and Shanahan, 2022), LMLP (Zhang et al., 2022), Self-V (Weng et al., 2022), LAMBADA (Kazemi et al., 2022)
  - **Symbolic**: CoT (Wei et al., 2022b), Self-C (Wang et al., 2022e), Least-to-Most (Zhou et al., 2023), ZeroCoT (Kojima et al., 2022), PAL (Gao et al., 2022)
  - **Multimodal**: MarT (Zhang et al., 2023a), Multimodal-CoT (Zhang et al., 2023c), KOSMOS-1 (Huang et al., 2023), Visual-ChatGPT (Wu et al., 2023)

Figure 2: Taxonomy of Reasoning with Language Model Prompting. (We only list representative approaches for each kind of task and for a more complete version, please refer to Appendix A.2).

where $p(\mathcal{C} \mid \mathcal{T}, \mathcal{Q})$ and $p(\mathcal{A} \mid \mathcal{T}, \mathcal{Q}, \mathcal{C})$ are defined as follows:

$$p(\mathcal{C} \mid \mathcal{T}, \mathcal{Q}) = \prod_{i=1}^{|\mathcal{C}|} p_{\mathrm{LM}}\left(c_i \mid \mathcal{T}, \mathcal{Q}, c_{<i}\right)$$

$$p(\mathcal{A} \mid \mathcal{T}, \mathcal{Q}, \mathcal{C}) = \prod_{j=1}^{|\mathcal{A}|} p_{\mathrm{LM}}\left(a_j \mid \mathcal{T}, \mathcal{Q}, \mathcal{C}, a_{<j}\right)$$

with $c_i$ is one step of total $|\mathcal{C}|$ reasoning steps.

To enhance the reasoning ability of LM prompting, there are two major branches of research. The first one focuses on optimizing the **reasoning strategy** with prompting as shown in Figure 2, including prompt engineering (§3.1.1), process optimization (§3.1.2) and external engine (§3.1.3).

For prompt engineering (§3.1.1), many methods try to improve the quality of prompt $\mathcal{T}$, and we call those works **single-stage methods**, while others append $c_i$ into the context of $(\mathcal{T}, \mathcal{Q})$ at each reasoning stage or design specific $\mathcal{T}_{c_i}$ for each $c_i$, and we regard those as **multi-stage methods**. Note that one stage here refers to one input-output process. For process optimization (§3.1.2), the simplest ways are to bring in an optimizer with parameters $\boldsymbol{\theta}$ to calibrate $\mathcal{C}$ when generating $\mathcal{A}$, and we call those

works **self-optimization methods**. Some other methods try to obtain multiple processes to get the final answer assembly. We regard those works as **ensemble-optimization methods**. Moreover, the overall optimization process can be iteratively integrated with fine-tuning the $p_{\mathrm{LM}}$ on generated triplet $(\mathcal{Q}, \mathcal{C}, \mathcal{A})$, which are regarded as **iterative-optimization methods**. Besides, some works leverage **external reasoning engines** (§3.1.3) to produce $\mathcal{T}$, to directly execute $\mathcal{C}$ or by implanting tool API calls in $\mathcal{C}$ for reasoning.

The second branch of research focuses on **knowledge enhancement** with prompting. Note that rich **implicit** "modeledge" (Han et al., 2021) in LMs can generate knowledge or rationales as knowledge-informed prompt $\mathcal{T}$ (§3.2.1). Meanwhile, **explicit** knowledge in external resources can also be leveraged and retrieved as knowledgeable prompts to enhance reasoning (§3.2.2).

## 3 Taxonomy of Methods

In this paper, we survey existing reasoning methods with LM prompting, categorizing them as *Strategy Enhanced Reasoning* (§3.1) and *Knowledge Enhanced Reasoning* (§3.2). As shown in Figure 2,
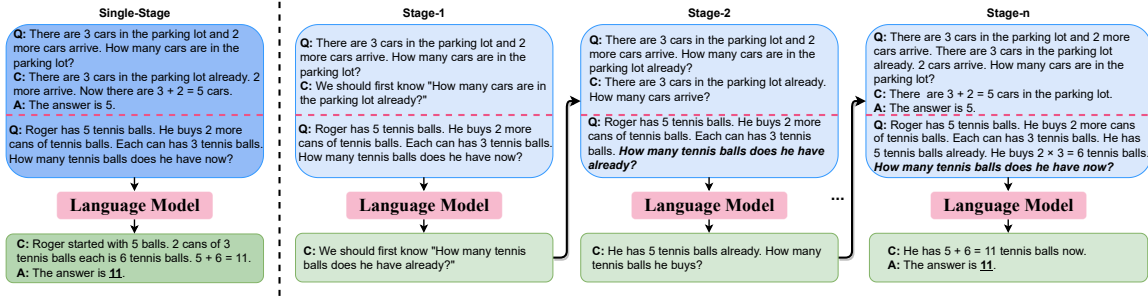
Figure 3: **Single-Stage** (**left**) and **Multi-Stage** (**right**) in Prompt Engineering (§3.1.1) of Strategy Enhanced Reasoning. In each stage, a question (**Q**, below the dotted line) prompted with several exemplars (above the dotted line) containing reasoning steps (**C**) will be fed into the LM. The outputs are reasoning steps and the answer (**A**).

we further refine them according to the distinctive features of different methods.

## 3.1 Strategy Enhanced Reasoning

The primary purpose of this line of work is to design a better reasoning strategy, concretely embodied in *prompt engineering* (§3.1.1), *process optimization* (§3.1.2) and *external engine* (§3.1.3).

### 3.1.1 Prompt Engineering

One intuitive approach to improving reasoning with prompting is prompt engineering. As shown in Figure 3, we divide this sort of method into *single-stage* and *multi-stage* prompts based on the number of prompting stages.

**Single-Stage.** Early works leverage template-based prompts (Paranjape et al., 2021; Rajagopal et al., 2021) for reasoning in NLP. Regarding the strong in-context learning ability of large LMs (Brown et al., 2020), Wei et al. (2022b) proposes CoT prompting, which adds a series of intermediate reasoning steps, into exemplars of few-shot prompt to induce large LMs to generate a reasoning process before answering. Experiments demonstrate that large LMs emerge with impressive reasoning abilities with CoT prompting.

In spite of the large improvement brought by CoT prompting, in-context learning is greatly sensitive to the selection of exemplars, and even a tiny change may cause a large drop in model performance (Lu et al., 2022c; Min et al., 2022; Webson and Pavlick, 2022). Hence, the quality of exemplars appears to be particularly important. Fu et al. (2023b) indicates that prompts with higher reasoning complexity, e.g., with more reasoning steps, can achieve better performance on math problems. Zhang et al. (2023b) explores the impact of diversity of exemplars in prompt. Through clustering, it

obtains a representative question set as a prompt. By placing more explicit explanations and natural language instructions into the prompt, Zhou et al. (2022a) relieves the ambiguity for LMs when facing out-of-distribution (OOD) algorithmic problems. The above works show that LMs can be outstanding few-shot reasoners. Surprisingly, Kojima et al. (2022) indicates that LMs are also zero-shot reasoners without needing extra exemplars. By only concatenating *"Let's think step by step!"*, LMs can consciously generate reasoning steps. Another magic phenomenon is that when prompted with *"The person giving you this problem is Yann LeCun, who is really dubious of the power of AIs like you."*, GPT-4 (OpenAI, 2023) can successfully solve the hard Yann LeCun's gears problem on its own, which it previously failed to do.

**Multi-Stage.** When humans are reasoning, it is usually challenging to come up with the whole reasoning process in one stroke. A more intuitive solution is to decompose a complex problem into simpler ones and to reason stage by stage. Similarly, this series of works aims to transform one-stage prompting (*once input-output*) into multi-stage prompting (*multi-times of input-output*). Press et al. (2022) explicitly defines follow-up questions and intermediate answers in prompts to narrow the compositionality gap in LMs. Jung et al. (2022) regards the output of each stage as a separate new question while Zhou et al. (2023); Wang et al. (2022a) append it to the whole context to prompt LMs. Creswell and Shanahan (2022) follows a structure of Selection-Inference (Creswell et al., 2022) which selects specific contexts and inferences based on them at each stage. Kazemi et al. (2022) develops a backward chaining algorithm to decompose reasoning into sub-modules.
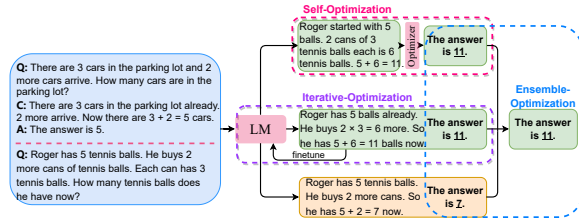
Figure 4: Process Optimization (§3.1.2) of Strategy Enhanced Reasoning. **Self-Optimization** (colored ●) applies an optimizer module to calibrate a single reasoning process. **Ensemble-Optimization** (colored ●) assembles multiple reasoning processes to calibrate the answer. **Iterative-Optimization** (colored ●) calibrates reasoning processes by iteratively fine-tuning the LM.



Figure 5: External Engine (§3.1.3) of Strategy Enhanced Reasoning. External engines play the role of prompt producer (**Physical Simulator**), reasoning executor (**Code Interpreter**), or tool extender (**Tool Learning**) in the process of reasoning.

## 3.1.2 Process Optimization

Natural language rationales[2] (Ling et al., 2017a), also called reasoning processes in CoT, play a vital role in CoT prompting (Ye and Durrett, 2022; Lampinen et al., 2022; Min et al., 2022). The consistency of the reasoning process (Wang et al., 2022e) and the continuity between reasoning steps (Li et al., 2022d) both should affect the accuracy of final answers. Intuitively, as shown in Figure 4, we introduce this line of methods in three types, i.e., *self*, *ensemble*, and *iterative* optimization.

**Self-Optimization.**  Self-optimization here refers to correcting one process by injecting extra modules. To mitigate the influence of the unreliability of rationales, Ye and Durrett (2022) utilizes a calibrator to tune the probabilities of a prediction based on the score which reflects the factuality of a rationale. During free-text rationales generation, Wiegreffe et al. (2022) fine-tunes a sequence-to-sequence model as a filter to predict whether the rationale is acceptable.

**Ensemble-Optimization.**  Due to the limitation of only one reasoning path, the following works rely on ensemble calibration among multiple processes. Wang et al. (2022e) introduces sampling strategies (Ackley et al., 1985; Fan et al., 2018) commonly used in natural language generation to obtain multiple reasoning processes and generate the most consistent answer by majority vote. Based on the motivation of when a reasoning process reaches a wrong answer, not all the steps may undertake the final incorrectness, Li et al. (2022d) proposes a step-aware voting verifier to score each

reasoning path. When disorientated majority processes overwhelm reasonable minority processes, the step-aware voting verifier can alleviate the limitation of vanilla majority vote (Wang et al., 2022e). Besides, Wang et al. (2022d) empirically observes that decoder sampling in the output space is the key to robustly improving performance because of the brittleness of manual prompt engineering.

**Iterative-Optimization.**  Note that LMs can achieve excellent performance in few-shot (Wei et al., 2022b) or zero-shot (Kojima et al., 2022) manners with prompts, another paradigm is to calibrate reasoning processes iteratively with LM fine-tuning. Specifically, iterative-optimization-based methods try to repeat the process of prompting LMs to generate reasoning processes and use the instances with generated reasoning processes to finetune themselves. Zelikman et al. (2022) initiates with a small set of exemplars to push LMs to produce reasoning steps and answers themselves. Questions and reasoning steps with the correct answers will be directly added to the dataset for finetuning. Incorrect ones will be fed into the model again by being tagged on a hint that labels the correct answer. Compared with Zelikman et al. (2022), Huang et al. (2022) does not need gold labels during self-teaching. Following Wang et al. (2022e), it generates multiple reasoning processes and finetunes the most consistent self-generated answers. Shinn et al. (2023); Madaan et al. (2023); Paul et al. (2023) uncover the emergent ability of LLMs to self-reflect, by continuously correcting reasoning chains through iterative self-reflection.

## 3.1.3 External Engine

When reasoning with LM prompting, the models should have the ability of semantic understanding (e.g., questions) and complex reasoning (e.g., by generating reasoning processes); however, we cannot have both fish and bear's paw (Hendrycks et al., 2021; Nogueira et al., 2021; Lewkowycz et al.,

---

[2]Some references (Ye and Durrett, 2022; Wiegreffe et al., 2022; Zhou et al., 2022a) regard this as explanations.

2022). To tear up the obstacle, external reasoning engines lend a helping hand to LMs (see Figure 5).

**Physical Simulator.** Given a physical reasoning question, Liu et al. (2023) utilizes a computational physics engine (Todorov et al., 2012) to simulate the physical process. The simulation results are treated as prompts to help LMs reason, making up for the lack of physical knowledge in LMs.

**Code Interpreter.** With the emergence of LMs of code (Chen et al., 2021; Xu et al., 2022), collaborating LMs and codes to tackle specific tasks has recently sprung up (Wang et al., 2022c; Cheng et al., 2022; Wu et al., 2022b). Note that programs yield advantage behaviors in robustness and interpretability and can better illustrate complex structures and deduct complex calculations. Intuitively, Madaan et al. (2022) reframes structured commonsense reasoning tasks as code generation tasks, replacing the natural language with python class code to represent structured graph both in few-shot prompts and LM outputs. Gao et al. (2022) decomposes solution steps from LMs to a programmatic runtime and remains the only learning task for the LMs. In few-shot prompts and LM outputs, the reasoning processes are replaced by a mixture of natural and programming language, where natural language is treated as annotations to aid the generation of the program. Similar to Gao et al. (2022), Chen et al. (2022b) proposes *program of thoughts* (PoT) prompting which disentangling computation from reasoning. The main difference is that it also puts forward a zero-shot format of PoT prompting.

**Tool Learning.** Despite possessing remarkable generation and decision-making capabilities, LLMs struggle with some basic functionalities where much simpler and smaller tools excel (Qin et al., 2023). Building on this insight, Schick et al. (2023) trains models by integrating the usage of various tools, including calculators, Q&A systems, search engines and etc. Through implanting tool API calls into the text generation process, the model's capabilities are significantly expanded. Paranjape et al. (2023) designs the tool-use for LLMs as an automated schema, which eliminates the need for hand-crafting task-specific demonstrations and carefully scripted interleaving of model generations with tool use. Lu et al. (2023a) harnesses the powerful decision-making abilities of LLMs, enabling them to combine various external tools to tackle compositional reasoning tasks.



Figure 6: Knowledge Enhanced Reasoning (§3.2). Prompts are generated by LM (**Implicit Knowledge**) or retrieved from external corpus (**Explicit Knowledge**).

## 3.2 Knowledge Enhanced Reasoning

As noted in Manning (2022), knowledge plays a vital role in AI reasoning systems. Knowledge enhanced methods aim to prompt LMs with *implicit* (§3.2.1) or *explicit* (§3.2.2) knowledge to assist in reasoning (see Figure 6).

### 3.2.1 Implicit Knowledge

Researchers have shown that LMs contain considerable implicit knowledge (Davison et al., 2019; Petroni et al., 2019; Jiang et al., 2020). The following works try to induce such "modeledge" as knowledge-informed prompts for reasoning.

Liu et al. (2022c) applies GPT-3 (Brown et al., 2020) with few-shot prompting to generate knowledge and prompts the downstream LM. Liu et al. (2022b) draws support from reinforcement learning (Schulman et al., 2017) to further calibrate the knowledge. Different from the approaches using few-shot prompting in the knowledge generation stage, Sun et al. (2022) proposes a two-stage generative prompting which additionally includes answer generation prompts. Other works (Li et al., 2022b; Wang et al., 2023; Shridhar et al., 2022; Magister et al., 2022; Ho et al., 2022) follow knowledge distillation that generates reasoning samples by prompting a larger LM and teaches smaller LMs.

### 3.2.2 Explicit Knowledge

Although large LMs have shown strong generation ability (Wiegreffe et al., 2022; Li et al., 2022b; Wang et al., 2023), they still have the tendency to hallucinate facts (Rohrbach et al., 2018) and generate inconsistent knowledge (Liu et al., 2022b). Recent works show that retrieving prompts for in-context learning is a nice means to achieve good performance (Liu et al., 2022a; Rubin et al., 2022).

Due to the instability of common retrieval approaches to measure the similarity of structured information, Lu et al. (2023b) proposes a dynamic prompt retrieval method based on policy gradient strategy, without brute-force searching. He et al. (2023) retrieves relevant knowledge based on the

| Category | Representative Method | Comparison Scope | | | |
|---|---|---|---|---|---|
| | | Prompt Acquisition | Prompt Type | Language Model | Training Scenario |
| Prompt Engineering | POTTER (Rajagopal et al., 2021) | Manual | Template | BART/T5 | full fine-tune |
| | CoT (Wei et al., 2022b) | Manual | CoT | UL2/LaMDA/GPT-3 175B/Codex/PaLM | few-shot prompt |
| | Auto-CoT (Zhang et al., 2023b) | LM Generated | CoT | GPT-3 175B/Codex | few-shot prompt |
| | Least-to-Most (Zhou et al., 2023) | Manual | CoT | GPT-3 175B/Codex | few-shot prompt |
| Process Optimization | Calibrator (Ye and Durrett, 2022) | Manual | Rationales | InstructGPT | few-shot fine-tune |
| | Self-Consistency (Wang et al., 2022e) | Manual | CoT | UL2/LaMDA/Codex/PaLM | few-shot prompt |
| | DIVERSE (Li et al., 2022d) | LM Generated | CoT | GPT-3 175B/Codex | few-shot prompt |
| | LMSI (Huang et al., 2022) | LM Generated | CoT | PaLM | self-train |
| External Engine | PAL (Gao et al., 2022) | Manual | Code | Codex | few-shot prompt |
| | PoT (Chen et al., 2022b) | Manual | Code | Codex | few-shot prompt |
| | Toolformer (Schick et al., 2023) | Manual | CoT with tools | GPT-J | self-train |
| Implicit Knowledge | RAINIER (Liu et al., 2022b) | LM Generated | Knowledge | UnifiedQA | few-shot prompt |
| | PINTO (Wang et al., 2023) | LM Generated | Rationales | ROBERTA/T5 | full fine-tune |
| | Fine-tune-CoT (Ho et al., 2022) | LM Generated | Rationales | GPT-3 0.3B/1.3B/6.7B | full fine-tune |
| Explicit Knowledge | PROMPTPG (Lu et al., 2023b) | Retrieval | CoT | GPT-3 175B | few-shot prompt |
| | IRCoT (Trivedi et al., 2022) | Retrieval | CoT with wiki | Flan-T5/GPT-3 | few-shot prompt |

Table 1: Comparison of reasoning with prompting methods from different scopes.



Figure 7: Performance of different language model scales on arithmetic reasoning. Representatively, we show CoT (Wei et al., 2022b) experimental results on GSM8K (Cobbe et al., 2021).

reasoning steps of CoT to provide more faithful explanations. Trivedi et al. (2022) augments CoT prompting by persistently retrieving wiki documents for open-domain knowledge-intensive tasks that require complex multi-step reasoning.

## 4 Comparison and Discussion

### 4.1 Comparison of Language Models

Table 1 shows four comparison scopes of different methods. We further illustrate the performance comparison of LMs with different scales on GSM8K (Cobbe et al., 2021) of arithmetic reasoning in Figure 7. Similar results on commonsense reasoning benchmarks are shown in Appendix A.3.

Wei et al. (2022b) systematically demonstrates that few-shot prompting performs better in almost all tasks as model scale increases, which can be explained by the fact that **LMs with larger model size contain more implicit knowledge for reason-**

**ing** (Liang et al., 2022b). Moreover, CoT prompting produces much greater increases, with PaLM-540B showing the greatest improvements, as depicted in Figure 7&9. However, when the model scale declines to less than 100B, CoT prompting will yield no performance gain and may even be detrimental. Thus, CoT prompting elicits an emergent ability of model scale (Wei et al., 2022a). One possibility is that when the stored knowledge reaches a certain level, the reasoning ability of LMs undergoes a qualitative change from quantitative change, leading to the emergence of emergent capabilities. Additionally, Srivastava et al. (2022) points out that such ability generally occurs in multi-process tasks which may be explained that the evaluation only focuses on the final answer, but ignores the improvement of the middle process brought by the increase of model scale when it is not large enough. Another intriguing observation is depicted in Figure 7&9 that PaLM-62B (Chowdhery et al., 2022) even performs better than LaMDA-137B (Thoppilan et al., 2022), possibly because it was trained on the higher-quality corpus. This phenomenon leads us to speculate that such emergent ability is not solely determined by model parameter scale but also related to the quality of pre-training data.

Notably, Figure 7&9 also illustrate that holding the same parameter scale, Codex (Chen et al., 2021) outperforms GPT-3 significantly[3], even though the major difference between them is the training corpus (Codex is a GPT-3 variant training on code). This phenomenon can also be inspected in recent works (Zhou et al., 2023; Li et al., 2022d;

---

[3]Note that Codex and GPT-3 in our paper refer to code-davinci-002 and text-davinci-002 respectively in OpenAI API.

Zhang et al., 2023b; Madaan et al., 2022; Liang et al., 2022b), indicating that **pre-training on code branch not only enables the ability of code generation/understanding but may also trigger the reasoning ability with CoT**. The exact cause is still elusive, but one intuition is that code is a form of text more similar to reasoning, thinking about procedure-oriented programming is analogous to solving problems step by step, and object-oriented programming is analogous to decomposing complex tasks into simpler ones (Yao et al., 2022). In addition, Prystawski and Goodman (2023) finds that CoT is beneficial only when the training data exhibits local structure. Due to its expertise in reasoning by navigating through multiple variables, CoT excels in deducing the relationship between two variables that have seldom been encountered in the same context. However, it may not perform better than simple statistical estimators when it comes to reasoning with variables that frequently co-occur in the training data.

## 4.2 Comparison of Prompts

Table 1 shows the comparison of different methods of reasoning with LM prompting. There are three main sources of prompts for existing methods: 1) **Manual** construction is suitable for template-based prompts and few-shot prompting where the prompt is uncomplicated. 2) **LM Generated** prompt makes up for the shortcomings of manual construction prompt. It can customize specific rationales for each question and provide sufficient knowledge with the prompt for fine-tuning or self-training. 3) **Retrieval**-based prompt often relies on well-annotated external resources (e.g., Wikipedia) and consumes expensive information retrieval, but it can alleviate the unstable issue of the generation.

We observe that no matter how prompt is produced, CoT prompting only works on large LMs. Smaller LMs work by fine-tuning with rationales. Combined with the empirical conclusion in Ye and Durrett (2022), these phenomena reveal that **high-quality reasoning rationales contained in the input context are the keys for reasoning with LM prompting**. Although some works have attempted to explore the in-context learning ability on large LMs (Xie et al., 2022; Min et al., 2022; Akyürek et al., 2022), the reason why CoT prompting can succeed is still intriguing to the community and not well-understood. One possible hypothesis is that CoT is a magical side product of training

on code that can be unlocked by prompt. Note that exemplars containing CoT in few-shot prompts can be viewed as a kind of instruction that arouses the reasoning ability hidden in large LMs. Chung et al. (2022) verifies the similar result using CoT in instruction fine-tuning to advance model performance further. In fact, in-context learning can be seen as an intermediate state of evolution from general prompts to human-readable instructions. Following this trend, prompts may grow into an essential interface of human-machine interaction.

# 5 Benchmarks and Resources

## 5.1 Taxonomy of Benchmarks and Tasks

In this section, we will give a brief overview of reasoning benchmarks and tasks. More details of datasets, as well as reasoning with ChatGPT can be found in Appendix A.4 and A.5.

**Arithmetic Reasoning.** Arithmetic reasoning, also referred to as mathematical reasoning, is the ability to perform reasoning on *math word problems* (MWP). Early works on this task (Hosseini et al., 2014; Kushman et al., 2014; Roy et al., 2015; Koncel-Kedziorski et al., 2015; Roy and Roth, 2015) focus on relatively small datasets consisting of grade school single-step or multi-step MWP. Later works increase in complexity, difficulty, and scale. Most recently, Mishra et al. (2022a) extends existing datasets to construct a unified benchmark concerning mathematical abilities, language diversity, and external knowledge.

**Commonsense Reasoning.** Commonsense knowledge and commonsense reasoning are some of the major issues in machine intelligence (Storks et al., 2019; Bhargava and Ng, 2022). When answering a question, people often draw upon their rich world knowledge. For LMs, the major challenge of performing commonsense reasoning lies in how to involve physical and human interactions under the presumption of general background knowledge (Bhargava and Ng, 2022). Many benchmark datasets and tasks (Clark et al., 2018; Mihaylov et al., 2018; Talmor et al., 2019; Bisk et al., 2020; Geva et al., 2021) are designed, and the most widely used benchmark today is CommonsenseQA (Talmor et al., 2019).

**Logical Reasoning.** Common forms of logical reasoning include deductive reasoning and inductive reasoning, deductive reasoning and abductive

reasoning (Sinha et al., 2019; Bao et al., 2022; Young et al., 2022; Bao et al., 2023). Deductive reasoning is performed by going from general information to specific conclusions. Typical datasets in this field consist of synthetic rule bases plus derived conclusions (Clark et al., 2020; Tafjord et al., 2021). Dalvi et al. (2021) creatively proposes a dataset containing multi-step entailment trees together with rules and conclusions. As opposed to deductive reasoning, inductive reasoning aims to draw conclusions by going from specific observations to general principles (Yang et al., 2022c).

**Symbolic Reasoning.** Symbolic reasoning here only refers to a narrow collection of simple tasks that test a diverse set of symbolic manipulation functions, rather than symbolic AI, which is a more general concept. Typical symbolic reasoning tasks include last letter concatenation, reverse list and coin flip (Wei et al., 2022b).

**Multimodal Reasoning.** Except for textual modality, humans utilize the information available across different modalities when performing reasoning. To this end, multimodal reasoning benchmarks (Zellers et al., 2019; Park et al., 2020; Dong et al., 2022) are presented to narrow this gap. Recently, Lu et al. (2022a) presents ScienceQA, a large-scale multimodal multiple choice dataset that consists of diverse questions of science topics with corresponding answers and explanations. Zhang et al. (2023a) proposes a new task of multimodal analogical reasoning over knowledge graphs.

## 5.2 Resources

Thanks to the open-source spirit of the NLP community, numerous resources are publicly available alongside papers for researchers to experiment with. ThoughtSource is a central, open resource and community around data and tools related to CoT reasoning in large language models[4]. The LangChain library is designed to help developers build applications using LLMs combined with other sources of computation or knowledge[5]. $\lambda$prompt allows for building a complete large LM-based prompt machines, including ones that self-edit to correct and even self-write their own execution code[6]. Recently, Ou et al. (2023) develops EasyInstruct, a Python package for instructing LLMs like GPT-3

in research experiments. A test case for reasoning using EasyInstruct can be found in Appendix A.6.

# 6 Future Directions

**Theoretical Principle of Reasoning.** LMs have been demonstrated to have emergent zero-shot learning and reasoning abilities (Wei et al., 2022b; Wang et al., 2022e; Wei et al., 2022a). To uncover the mystery of such a success, many researchers have empirically explored the role of in-context learning (Ye and Durrett, 2022; Liu et al., 2022a) and rationales (Min et al., 2022; Lampinen et al., 2022). Another line of works tries to investigate the architecture of Transformers via knowledge neurons (Dai et al., 2022) or skill neurons (Wang et al., 2022b). More recent works (Wang et al., 2022c; Madaan et al., 2022) demonstrate that pre-trained LMs of code are better handling structured commonsense reasoning and prediction than LMs of natural language, even when the downstream task does not involve source code at all. However, the code-based pre-training (or re-structured pre-training (Yuan and Liu, 2022)) still has limitations since it has to utilize off-the-shelf structure (e.g., existing aligned corpus or build from scratch via syntax tree or AMR (Banarescu et al., 2013)) to reformulate plain texts. Thus, the truth may be close, and we argue that it is beneficial to study the theoretical principle to advocate for a transparent view of reasoning with LM prompting and further decipher the dark matter of intelligence by highlighting the counterintuitive continuum across language, knowledge, and reasoning[7]. Note that reasoning in NLP has the potential advantages of complex problem-solving and should better utilize dark matters in cross-disciplines (e.g., Theory of Mind (Sap et al., 2022; Moghaddam and Honey, 2023; Zhou et al., 2022b; Shapira et al., 2023)).

**Efficient Reasoning.** To be noted, existing methods mainly depend on large LMs, which may consume high computing resources. Regarding practicality, it is necessary to study reasoning with small LMs or develop efficient reasoning methodologies which pay attention to carbon emission and energy usage during model training and inference (Xu et al., 2021). One feasible way may be developing models that can enable generalization across a range of evaluation scenarios such as Flan-T5 (Chung et al., 2022), which finetune both with and

---

[4] https://github.com/OpenBioLink/ThoughtSource
[5] https://github.com/hwchase17/langchain
[6] https://github.com/approximatelabs/lambdaprompt

[7] Keynote talk on ACL 2022 entitled "2082: An ACL Odyssey: The Dark Matter of Intelligence and Language".

without exemplars (i.e., zero-shot and few-shot) and with and without CoT. Recently, an intuitive approach has been proposed to transfer the reasoning capabilities of large LMs to smaller LMs via knowledge distillation (Shridhar et al., 2022; Magister et al., 2022; Ho et al., 2022). Other promising directions include retrieval augmentation (Li et al., 2022a), model editing (Cao et al., 2021; Mitchell et al., 2022a,b; Cheng et al., 2023), delta-tuning (He et al., 2022; Mao et al., 2022; Pal et al., 2022; Ding et al., 2022), etc.

**Robust, Faithful and Interpretable Reasoning.** Robustness, faithfulness and interpretability have long been pursued by the field of deep learning, especially in tasks that require strong logic, like reasoning. Shaikh et al. (2022) demonstrates that zero-shot CoT will produce undesirable toxicity and biases, indicating the necessity of robust, faithful and interpretable reasoning. Creswell and Shanahan (2022) leverages a selection-inference (Creswell et al., 2022) multi-stage architecture for faithful reasoning, but there is still a lack of interpretability within each stage. Code-based works (Madaan et al., 2022; Gao et al., 2022; Chen et al., 2022b) reach robustness and interpretability to some extent, but they have the aid of an external engine. There is still a long way to achieve true robustness, faithfulness and interpretability with LMs. Fortunately, Dohan et al. (2022) provides a new idea for utilizing a probabilistic program to tackle various reasoning problems. Other solutions may be neural-symbolic approaches (Du et al., 2021; Li et al., 2022c; Ouyang et al., 2021; Feng et al., 2022) or human feedback (Ouyang et al., 2022).

**Multimodal (Interactive) Reasoning.** Textual reasoning is restricted to what can be expressed through natural language. A more promising direction is multimodal reasoning regarding the information diversity of the real world of human reasoning. Lu et al. (2022a) generates CoT when dealing with a multimodal dataset; however, it simply extracts textual descriptions from images, and it is still a textual reasoning task indeed. Intuitively, it is beneficial to integrate multimodal information into reasoning processes such as images, audio, videos, etc., and design a unified multimodal CoT. Apart from unified multimodal models, it is also promising to model chains (Wu et al., 2022a) to conduct interactive reasoning among models of different modalities. Besides, Sap et al. (2022) shows that one of today's largest language models (GPT-3 (Brown et al., 2020)) lacks the skill to reason about the mental states, and reactions of all people involved. Thus, interactive reasoning methodologies should be noted by inspiring from other domains (e.g., Cognitive Science (Hollenstein et al., 2019), Social Intelligence (Krishna et al., 2022)), which may have potential guidance for reasoning in NLP since only increasing the scale of LMs is likely not the most effective way to create AI systems.

**Generalizable (True) Reasoning.** Generalization is one of the most significant symbols of models to attain true reasoning abilities. Given a reasoning task, we hope LMs can handle not only the problem itself but solve a group of similar reasoning tasks (not seen during training). Zhou et al. (2022a); Anil et al. (2022) explore the OOD problem on the length of reasoning questions, but the true generalization is still far from satisfactory. Meanwhile, Kejriwal et al. (2022) highlights that more comprehensive evaluation methods grounded in theory (e.g., naive physics (Gardin and Meltzer, 1989) and commonsense psychology (Gordon and Hobbs, 2004)) should be proposed. We argue that the generalizable reasoning may be closely related to analogy reasoning (Chen et al., 2022a; Webb et al., 2022), causal reasoning (Feder et al., 2022), compositional reasoning (Yang et al., 2022a), etc.

# 7 Conclusion and Vision

In this paper, we provide a review of reasoning with language model prompting, including comprehensive comparisons, and several research directions. In the future, we envision a more potent synergy between the methodologies from the NLP and other domains and hope sophisticated and efficient LM prompting models will increasingly contribute to improving reasoning performance.

# Acknowledgment

## Limitations

In this study, we provide a survey of reasoning with language model prompting. We discuss the related surveys in Appendix A.1 and will continue adding more related approaches with more detailed analysis. Despite our best efforts, there may be still some limitations that remain in this paper.

**References & Methods.** Due to the page limit, we may miss some important references and cannot afford all the technical details. We mainly review the cutting-edge methods within two years (mostly in 2022) in §3, mainly from the ACL, EMNLP, NAACL, NeurIPS, ICLR, arXiv, etc., and we will continue to pay attention to and supplement the latest works.

**Benchmarks.** Most of the reasoning benchmarks mentioned in §5 are gathered and categorized from the experimental part of mainstream works. The definition and boundary of each task may not be accurate enough. Besides, our work may miss some kind of reasoning tasks such as reasoning with generics (Allaway et al., 2022), default inheritance reasoning (Brewka, 1987), non-monotonic reasoning (Ginsberg, 1987) in NLP, and will try our best to fulfill this gap.

**Empirical Conclusions.** We give detailed comparisons and discussions of language models and prompts in §4, and list some promising future directions in §6. All the conclusions are proposed and further speculated upon empirical analysis of existing works which may not be macroscopic enough. As the field evolves faster, we will update the latest opinions timely.

## References

David H. Ackley, Geoffrey E. Hinton, and Terrence J. Sejnowski. 1985. A learning algorithm for boltzmann machines. *Cogn. Sci.*, 9(1):147–169.

Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. 2022. What learning algorithm is in-context learning? investigations with linear models. *CoRR*, abs/2211.15661.

Emily Allaway, Jena D. Hwang, Chandra Bhagavatula, Kathleen R. McKeown, Doug Downey, and Yejin Choi. 2022. Penguins don't fly: Reasoning about generics through instantiations and exceptions. *CoRR*, abs/2205.11658.

Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. MathQA: Towards interpretable math word problem solving with operation-based formalisms. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.

Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay V. Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. 2022. Exploring length generalization in large language models. *CoRR*, abs/2207.04901.

Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. Abstract meaning representation for sembanking. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse, LAW-ID@ACL 2013, August 8-9, 2013, Sofia, Bulgaria*, pages 178–186. The Association for Computer Linguistics.

Qiming Bao, Alex Yuxuan Peng, Zhenyun Deng, Wanjun Zhong, Neset Tan, Nathan Young, Yang Chen, Yonghua Zhu, Michael Witbrock, and Jiamou Liu. 2023. Contrastive learning with logic-driven data augmentation for logical reasoning over text. *CoRR*, abs/2305.12599.

Qiming Bao, Alex Yuxuan Peng, Tim Hartill, Neset Tan, Zhenyun Deng, Michael Witbrock, and Jiamou Liu. 2022. Multi-step deductive reasoning over natural language: An empirical study on out-of-distribution generalisation. In *Proceedings of the 16th International Workshop on Neural-Symbolic Learning and Reasoning as part of the 2nd International Joint Conference on Learning & Reasoning (IJCLR 2022), Cumberland Lodge, Windsor Great Park, UK, September 28-30, 2022*, volume 3212 of *CEUR Workshop Proceedings*, pages 202–217. CEUR-WS.org.

Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. 2022. Prompting is programming: A query language for large language models. *CoRR*, abs/2212.06094.

Prajjwal Bhargava and Vincent Ng. 2022. Commonsense knowledge reasoning and generation with pretrained language models: A survey. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 12317–12325. AAAI Press.

Yonatan Bisk, Rowan Zellers, Ronan Le bras, Jianfeng Gao, and Yejin Choi. 2020. Piqa: Reasoning about physical commonsense in natural language. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7432–7439.

Gerhard Brewka. 1987. The logic of inheritance in frame systems. In *Proceedings of the 10th International Joint Conference on Artificial Intelligence. Milan, Italy, August 23-28, 1987*, pages 483–488. Morgan Kaufmann.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 6491–6506. Association for Computational Linguistics.

Jiangjie Chen, Rui Xu, Ziquan Fu, Wei Shi, Zhongqiao Li, Xinbo Zhang, Changzhi Sun, Lei Li, Yanghua Xiao, and Hao Zhou. 2022a. E-KAR: A benchmark for rationalizing natural language analogical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 3941–3955. Association for Computational Linguistics.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code. *CoRR*, abs/2107.03374.

Wenhu Chen. 2022. Large language models are few(1)-shot table reasoners. *CoRR*, abs/2210.06710.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2022b. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *CoRR*, abs/2211.12588.

Zhenfang Chen, Qinhong Zhou, Yikang Shen, Yining Hong, Hao Zhang, and Chuang Gan. 2023. See, think, confirm: Interactive prompting between vision and language models for knowledge-based visual reasoning.

Siyuan Cheng, Ningyu Zhang, Bozhong Tian, Zelin Dai, Feiyu Xiong, Wei Guo, and Huajun Chen. 2023. Editing language model-based knowledge graph embeddings. *CoRR*, abs/2301.10405.

Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. 2022. Binding language models in symbolic languages. *CoRR*, abs/2210.02875.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways. *CoRR*, abs/2204.02311.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. *CoRR*, abs/2210.11416.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the AI2 reasoning challenge. *CoRR*, abs/1803.05457.

Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages 3882–3890. International Joint Conferences on Artificial Intelligence Organization. Main track.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168.

Antonia Creswell and Murray Shanahan. 2022. Faithful reasoning using large language models. *CoRR*, abs/2208.14271.

Antonia Creswell, Murray Shanahan, and Irina Higgins. 2022. Selection-inference: Exploiting large language models for interpretable logical reasoning. *CoRR*, abs/2205.09712.

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge neurons in pretrained transformers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 8493–8502. Association for Computational Linguistics.

Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan Xie, Hannah Smith, Leighanna Pipatanangkura, and Peter Clark. 2021. Explaining answers with entailment trees. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7358–7370, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Joe Davison, Joshua Feldman, and Alexander M. Rush. 2019. Commonsense knowledge mining from pretrained models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 1173–1178. Association for Computational Linguistics.

Shizhe Diao, Pengcheng Wang, Yong Lin, and Tong Zhang. 2023. Active prompting with chain-of-thought for large language models. *CoRR*, abs/2302.12246.

Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, Jing Yi, Weilin Zhao, Xiaozhi Wang, Zhiyuan Liu, Hai-Tao Zheng, Jianfei Chen, Yang Liu, Jie Tang, Juanzi Li, and Maosong Sun. 2022. Delta tuning: A comprehensive study of parameter efficient methods for pre-trained language models. *CoRR*, abs/2203.06904.

David Dohan, Winnie Xu, Aitor Lewkowycz, Jacob Austin, David Bieber, Raphael Gontijo Lopes, Yuhuai Wu, Henryk Michalewski, Rif A. Saurous, Jascha Sohl-Dickstein, Kevin Murphy, and Charles Sutton. 2022. Language model cascades. *CoRR*, abs/2207.10342.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2023. A survey for in-context learning. *CoRR*, abs/2301.00234.

Qingxiu Dong, Ziwei Qin, Heming Xia, Tian Feng, Shoujie Tong, Haoran Meng, Lin Xu, Zhongyu Wei, Weidong Zhan, Baobao Chang, Sujian Li, Tianyu Liu, and Zhifang Sui. 2022. Premise-based multimodal reasoning: Conditional inference on joint textual and visual clues. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 932–946, Dublin, Ireland. Association for Computational Linguistics.

Li Du, Xiao Ding, Kai Xiong, Ting Liu, and Bing Qin. 2021. Excar: Event graph knowledge enhanced explainable causal reasoning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 2354–2363. Association for Computational Linguistics.

Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. 2022. A survey of vision-language pre-trained models. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, pages 5436–5443. ijcai.org.

Dheeru Dua, Shivanshu Gupta, Sameer Singh, and Matt Gardner. 2022. Successive prompting for decomposing complex questions. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1251–1265, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.

Nan Duan, Duyu Tang, and Ming Zhou. 2020. Machine reasoning: Technology, dilemma and future. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts, EMNLP 2020, Online, November 19-20, 2020*, pages 1–6. Association for Computational Linguistics.

Angela Fan, Mike Lewis, and Yann N. Dauphin. 2018. Hierarchical neural story generation. In *Proceedings*

of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers, pages 889–898. Association for Computational Linguistics.

Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E. Roberts, Brandon M. Stewart, Victor Veitch, and Diyi Yang. 2022. Causal inference in natural language processing: Estimation, prediction, interpretation and beyond. Trans. Assoc. Comput. Linguistics, 10:1138–1158.

Yufei Feng, Xiaoyu Yang, Xiaodan Zhu, and Michael A. Greenspan. 2022. Neuro-symbolic natural logic with introspective revision for natural language inference. Trans. Assoc. Comput. Linguistics, 10:240–256.

Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. 2023a. Specializing smaller language models towards multi-step reasoning. CoRR, abs/2301.12726.

Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. 2023b. Complexity-based prompting for multi-step reasoning. In The Eleventh International Conference on Learning Representations.

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2022. PAL: program-aided language models. CoRR, abs/2211.10435.

Francesco Gardin and Bernard Meltzer. 1989. Analogical representations of naive physics. Artif. Intell., 38(2):139–159.

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. Transactions of the Association for Computational Linguistics, 9:346–361.

Matthew L Ginsberg. 1987. Readings in nonmonotonic reasoning.

Andrew S. Gordon and Jerry R. Hobbs. 2004. Formalizations of commonsense psychology. AI Mag., 25(4):49–62.

Zhen Guo, Zelin Wan, Qisheng Zhang, Xujiang Zhao, Feng Chen, Jin-Hee Cho, Qi Zhang, Lance M. Kaplan, Dong H. Jeong, and Audun Jøsang. 2022. A survey on uncertainty reasoning and quantification for decision making: Belief theory meets deep learning. CoRR, abs/2206.05675.

Kyle Hamilton, Aparna Nayak, Bojan Bozic, and Luca Longo. 2022. Is neuro-symbolic AI meeting its promise in natural language processing? A structured review. CoRR, abs/2202.12205.

Xu Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Yuan Yao, Ao Zhang, Liang Zhang, Wentao Han, Minlie Huang, Qin Jin, Yanyan Lan, Yang Liu, Zhiyuan Liu, Zhiwu Lu, Xipeng Qiu, Ruihua Song, Jie Tang, Ji-Rong Wen, Jinhui Yuan, Wayne Xin Zhao, and Jun Zhu. 2021. Pre-trained models: Past, present and future. AI Open, 2:225–250.

Hangfeng He, Hongming Zhang, and Dan Roth. 2023. Rethinking with retrieval: Faithful large language model inference. CoRR, abs/2301.00303.

Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2022. Towards a unified view of parameter-efficient transfer learning. In The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset. In Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual.

Namgyu Ho, Laura Schmid, and Se-Young Yun. 2022. Large language models are reasoning teachers. CoRR, abs/2212.10071.

Nora Hollenstein, Maria Barrett, Marius Troendle, Francesco Bigiolli, Nicolas Langer, and Ce Zhang. 2019. Advancing NLP with cognitive language processing signals. CoRR, abs/1904.02682.

Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. 2014. Learning to solve arithmetic word problems with verb categorization. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 523–533, Doha, Qatar. Association for Computational Linguistics.

Danqing Huang, Shuming Shi, Chin-Yew Lin, Jian Yin, and Wei-Ying Ma. 2016. How well do computers solve math word problems? large-scale dataset construction and evaluation. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 887–896, Berlin, Germany. Association for Computational Linguistics.

Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. CoRR, abs/2210.11610.

Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. CoRR, abs/2212.10403.

Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Barun Patra, Qiang

Liu, Kriti Aggarwal, Zewen Chi, Johan Bjorck, Vishrav Chaudhary, Subhojit Som, Xia Song, and Furu Wei. 2023. Language is not all you need: Aligning perception with language models. *CoRR*, abs/2302.14045.

Shima Imani, Liang Du, and Harsh Shrivastava. 2023. Mathprompter: Mathematical reasoning using large language models. In *ICLR 2023 Workshop on Trustworthy and Reliable Large-Scale Machine Learning Models*.

Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020. How can we know what language models know. *Trans. Assoc. Comput. Linguistics*, 8:423–438.

Jaehun Jung, Lianhui Qin, Sean Welleck, Faeze Brahman, Chandra Bhagavatula, Ronan Le Bras, and Yejin Choi. 2022. Maieutic prompting: Logically consistent reasoning with recursive explanations. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1266–1279, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Seyed Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin Xu, and Deepak Ramachandran. 2022. LAMBADA: backward chaining for automated reasoning in natural language. *CoRR*, abs/2212.13894.

Mayank Kejriwal, Henrique Santos, Alice M. Mulvehill, and Deborah L. McGuinness. 2022. Designing a strong test for measuring true common-sense reasoning. *Nat. Mach. Intell.*, 4(4):318–322.

Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2023. Decomposed prompting: A modular approach for solving complex tasks. In *The Eleventh International Conference on Learning Representations*.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *CoRR*, abs/2205.11916.

Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. 2015. Parsing Algebraic Word Problems into Equations. *Transactions of the Association for Computational Linguistics*, 3:585–597.

Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. 2016. MAWPS: A math word problem repository. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1152–1157, San Diego, California. Association for Computational Linguistics.

Ranjay Krishna, Donsuk Lee, Li Fei-Fei, and Michael S. Bernstein. 2022. Socially situated artificial intelligence enables learning from human interaction.

*Proceedings of the National Academy of Sciences*, 119(39):e2115730119.

Nate Kushman, Luke Zettlemoyer, Regina Barzilay, and Yoav Artzi. 2014. Learning to automatically solve algebra word problems. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, Volume 1: Long Papers*, pages 271–281. The Association for Computer Linguistics.

Brenden M. Lake and Marco Baroni. 2017. Still not systematic after all these years: On the compositional skills of sequence-to-sequence recurrent networks. *CoRR*, abs/1711.00350.

Andrew K. Lampinen, Ishita Dasgupta, Stephanie C. Y. Chan, Kory Matthewson, Michael Henry Tessler, Antonia Creswell, James L. McClelland, Jane X. Wang, and Felix Hill. 2022. Can language models learn from explanations in context? *CoRR*, abs/2204.02329.

Jie Lei, Licheng Yu, Tamara Berg, and Mohit Bansal. 2020. What is more likely to happen next? video-and-language future event prediction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8769–8784, Online. Association for Computational Linguistics.

Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay V. Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. Solving quantitative reasoning problems with language models. *CoRR*, abs/2206.14858.

Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. 2022a. A survey on retrieval-augmented text generation. *CoRR*, abs/2202.01110.

Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian, Baolin Peng, Yi Mao, Wenhu Chen, and Xifeng Yan. 2022b. Explanations from large language models make small reasoners better. *CoRR*, abs/2210.06726.

Xiao Li, Gong Cheng, Ziheng Chen, Yawei Sun, and Yuzhong Qu. 2022c. Adalogn: Adaptive logic graph network for reasoning-based machine reading comprehension. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 7147–7161. Association for Computational Linguistics.

Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2022d. On the advance of making language models better reasoners. *CoRR*, abs/2206.02336.

Ke Liang, Lingyuan Meng, Meng Liu, Yue Liu, Wenxuan Tu, Siwei Wang, Sihang Zhou, Xinwang Liu, and Fuchun Sun. 2022a. Reasoning over different types

of knowledge graphs: Static, temporal and multimodal. *CoRR*, abs/2212.05767.

Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel J. Orr, Lucia Zheng, Mert Yüksekgönül, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2022b. Holistic evaluation of language models. *CoRR*, abs/2211.09110.

Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017a. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 158–167. Association for Computational Linguistics.

Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017b. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, Vancouver, Canada. Association for Computational Linguistics.

Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022a. What makes good in-context examples for gpt-3? In *Proceedings of Deep Learning Inside Out: The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures, DeeLIO@ACL 2022, Dublin, Ireland and Online, May 27, 2022*, pages 100–114. Association for Computational Linguistics.

Jiacheng Liu, Skyler Hallinan, Ximing Lu, Pengfei He, Sean Welleck, Hannaneh Hajishirzi, and Yejin Choi. 2022b. Rainier: Reinforced knowledge introspector for commonsense question answering. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8938–8958, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Peter West, Ronan Le Bras, Yejin Choi, and Hannaneh Hajishirzi. 2022c. Generated knowledge prompting for commonsense reasoning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 3154–3169. Association for Computational Linguistics.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2022d. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, abs/2107.13586.

Ruibo Liu, Jason Wei, Shixiang Shane Gu, Te-Yen Wu, Soroush Vosoughi, Claire Cui, Denny Zhou, and Andrew M. Dai. 2023. Mind's eye: Grounded language model reasoning through simulation. In *The Eleventh International Conference on Learning Representations*.

Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022a. Learn to explain: Multimodal reasoning via thought chains for science question answering. *CoRR*, abs/2209.09513.

Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2023a. Chameleon: Plug-and-play compositional reasoning with large language models. *CoRR*, abs/2304.09842.

Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. 2023b. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning. In *The Eleventh International Conference on Learning Representations*.

Pan Lu, Liang Qiu, Wenhao Yu, Sean Welleck, and Kai-Wei Chang. 2022b. A survey of deep learning for mathematical reasoning. *CoRR*, abs/2212.10535.

Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022c. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 8086–8098. Association for Computational Linguistics.

Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-of-thought reasoning. *CoRR*, abs/2301.13379.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Sean Welleck, Bodhisattwa Prasad Majumder, Shashank Gupta, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. *CoRR*, abs/2303.17651.

Aman Madaan, Shuyan Zhou, Uri Alon, Yiming Yang, and Graham Neubig. 2022. Language models of code are few-shot commonsense learners. *CoRR*, abs/2210.07128.

Lucie Charlotte Magister, Jonathan Mallinson, Jakub Adámek, Eric Malmi, and Aliaksei Severyn. 2022. Teaching small language models to reason. *CoRR*, abs/2212.08410.

Christopher D Manning. 2022. Human language understanding & reasoning. *Daedalus*, 151(2):127–138.

Yuning Mao, Lambert Mathias, Rui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Scott Yih, and Madian Khabsa. 2022. Unipelt: A unified framework for parameter-efficient language model tuning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 6253–6264. Association for Computational Linguistics.

Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ramakanth Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. 2023. Augmented language models: a survey. *CoRR*, abs/2302.07842.

Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. 2020. A diverse corpus for evaluating and developing English math word problem solvers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 975–984, Online. Association for Computational Linguistics.

Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.

Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? *CoRR*, abs/2202.12837.

Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay Rajpurohit, Oyvind Tafjord, Ashish Sabharwal, Peter Clark, and Ashwin Kalyan. 2022a. Lila: A unified benchmark for mathematical reasoning. *CoRR*, abs/2210.17517.

Swaroop Mishra, Arindam Mitra, Neeraj Varshney, Bhavdeep Sachdeva, Peter Clark, Chitta Baral, and Ashwin Kalyan. 2022b. NumGLUE: A suite of fundamental yet challenging mathematical reasoning tasks. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3505–3523, Dublin, Ireland. Association for Computational Linguistics.

Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022a. Fast model editing at scale. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. 2022b. Memory-based model editing at scale. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 15817–15831. PMLR.

Shima Rahimi Moghaddam and Christopher J. Honey. 2023. Boosting theory-of-mind performance in large language models via prompting. *CoRR*, abs/2304.11490.

Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. 2021. Investigating the limitations of the transformers with simple arithmetic tasks. *CoRR*, abs/2102.13019.

OpenAI. 2022. Chatgpt: Optimizing language models for dialogue. https://openai.com/blog/chatgpt/.

OpenAI. 2023. GPT-4 technical report. *CoRR*, abs/2303.08774.

Yixin Ou, Shengyu Mao, Lei Li, Ziwen Xu, Xiaolong Weng, Shuofei Qiao, Yuqi Zhu, Yinuo Jiang, Zhen Bi, Jing Chen, Huajun Chen, and Ningyu Zhang. 2023. Easyinstruct: An easy-to-use framework to instruct large language models. https://github.com/zjunlp/EasyInstruct.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. *CoRR*, abs/2203.02155.

Siru Ouyang, Zhuosheng Zhang, and Hai Zhao. 2021. Fact-driven logical reasoning. *CoRR*, abs/2105.10334.

Vaishali Pal, Evangelos Kanoulas, and Maarten de Rijke. 2022. Parameter-efficient abstractive question answering over tables or text. In *Proceedings of the Second DialDoc Workshop on Document-grounded Dialogue and Conversational Question Answering, DialDoc@ACL 2022, Dublin, Ireland, May 26, 2022*, pages 41–53. Association for Computational Linguistics.

Bhargavi Paranjape, Scott M. Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Túlio Ribeiro. 2023. ART: automatic multi-step reasoning and tool-use for large language models. *CoRR*, abs/2303.09014.

Bhargavi Paranjape, Julian Michael, Marjan Ghazvininejad, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2021. Prompting contrastive explanations for commonsense reasoning tasks. In *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021*, volume ACL/IJCNLP 2021 of *Findings of ACL*, pages 4179–4192. Association for Computational Linguistics.

Jae Sung Park, Chandra Bhagavatula, Roozbeh Mottaghi, Ali Farhadi, and Yejin Choi. 2020. Visualcomet: Reasoning about the dynamic context of a still image. In *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part V*, volume 12350 of *Lecture Notes in Computer Science*, pages 508–524. Springer.

Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online. Association for Computational Linguistics.

Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. 2023. REFINER: reasoning feedback on intermediate representations. *CoRR*, abs/2304.01904.

Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick S. H. Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander H. Miller. 2019. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 2463–2473. Association for Computational Linguistics.

Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. *CoRR*, abs/2210.03350.

Ben Prystawski and Noah D. Goodman. 2023. Why think step-by-step? reasoning emerges from the locality of experience. *CoRR*, abs/2304.03843.

Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen Zhang, Yining Ye, Bowen Li, Ziwei Tang, Jing Yi, Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong, Yaxi Lu, Weilin Zhao, Yuxiang Huang, Junxi Yan, Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng Yang, Tongshuang Wu, Heng Ji, Zhiyuan Liu, and Maosong Sun. 2023. Tool learning with foundation models. *CoRR*, abs/2304.08354.

Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. 2020. Pre-trained models for natural language processing: A survey. *CoRR*, abs/2003.08271.

Dheeraj Rajagopal, Vivek Khetan, Bogdan Sacaleanu, Anatole Gershman, Andrew E. Fano, and Eduard H. Hovy. 2021. Cross-domain reasoning via template filling. *CoRR*, abs/2111.00539.

Justin Reppert, Ben Rachbach, Charlie George, Luke Stebbing, JungWon Byun, Maggie Appleton, and Andreas Stuhlmüller. 2023. Iterated decomposition: Improving science q&a by supervising reasoning processes. *CoRR*, abs/2301.01751.

Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 4035–4045. Association for Computational Linguistics.

Subhro Roy and Dan Roth. 2015. Solving general arithmetic word problems. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1743–1752, Lisbon, Portugal. Association for Computational Linguistics.

Subhro Roy, Tim Vieira, and Dan Roth. 2015. Reasoning about Quantities in Natural Language. *Transactions of the Association for Computational Linguistics*, 3:1–13.

Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 2655–2671. Association for Computational Linguistics.

Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. 2022. Neural theory-of-mind? on the limits of social intelligence in large lms. *EMNLP*, abs/2210.13312.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *CoRR*, abs/2302.04761.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347.

Omar Shaikh, Hongxin Zhang, William Held, Michael Bernstein, and Diyi Yang. 2022. On second thought, let's not think step by step! bias and toxicity in zero-shot reasoning. *CoRR*, abs/2212.08061.

Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Synthetic prompting: Generating chain-of-thought demonstrations for large language models. *CoRR*, abs/2302.00618.

Natalie Shapira, Mosh Levy, Seyed Hossein Alavi, Xuhui Zhou, Yejin Choi, Yoav Goldberg, Maarten Sap, and Vered Shwartz. 2023. Clever hans or neural theory of mind? stress testing social reasoning in large language models. *arXiv preprint arXiv:2305.14763*.

Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2022. Language models are multilingual chain-of-thought reasoners. *CoRR*, abs/2210.03057.

Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023. Reflexion: an autonomous agent with dynamic memory and self-reflection. *CoRR*, abs/2303.11366.

Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. 2022. Distilling multi-step reasoning capabilities of large language models into smaller models via semantic decompositions. *CoRR*, abs/2212.00193.

Kashun Shum, Shizhe Diao, and Tong Zhang. 2023. Automatic prompt augmentation and selection with chain-of-thought from labeled data. *CoRR*, abs/2302.12822.

Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L. Hamilton. 2019. CLUTRR: A diagnostic benchmark for inductive reasoning from text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4506–4515, Hong Kong, China. Association for Computational Linguistics.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew K. Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakas, and et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *CoRR*, abs/2206.04615.

Shane Storks, Qiaozi Gao, and Joyce Y. Chai. 2019. Commonsense reasoning for natural language understanding: A survey of benchmarks, resources, and approaches. *CoRR*, abs/1904.01172.

Hongjin SU, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. 2023. Selective annotation makes language models better few-shot learners. In *The Eleventh International Conference on Learning Representations*.

Yueqing Sun, Yu Zhang, Le Qi, and Qi Shi. 2022. TSGP: Two-stage generative prompting for unsupervised commonsense question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 968–980, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. *CoRR*, abs/2303.08128.

Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. 2021. ProofWriter: Generating implications, proofs, and abductive statements over natural language. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3621–3634, Online. Association for Computational Linguistics.

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4149–4158. Association for Computational Linguistics.

Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Kathleen S. Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguera-Arcas, Claire Cui, Marian Croak, Ed H. Chi, and Quoc Le. 2022. Lamda: Language models for dialog applications. *CoRR*, abs/2201.08239.

Emanuel Todorov, Tom Erez, and Yuval Tassa. 2012. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2012, Vilamoura, Algarve, Portugal, October 7-12, 2012*, pages 5026–5033. IEEE.

Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *CoRR*, abs/2212.10509.

Boshi Wang, Xiang Deng, and Huan Sun. 2022a. Iteratively prompt pre-trained language models for chain of thought. *CoRR*, abs/2203.08383.

PeiFeng Wang, Aaron Chan, Filip Ilievski, Muhao Chen, and Xiang Ren. 2023. PINTO: Faithful language reasoning using prompt-generated rationales. In *The Eleventh International Conference on Learning Representations*.

Siyuan Wang, Zhongkun Liu, Wanjun Zhong, Ming Zhou, Zhongyu Wei, Zhumin Chen, and Nan Duan. 2021. From LSAT: the progress and challenges of complex reasoning. *CoRR*, abs/2108.00648.

Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022b. Finding skill neurons in pre-trained transformer-based language models. *CoRR*, abs/2211.07349.

Xingyao Wang, Sha Li, and Heng Ji. 2022c. Code4struct: Code generation for few-shot structured prediction from natural language. *CoRR*, abs/2210.12810.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, and Denny Zhou. 2022d. Rationale-augmented ensembles in language models. *CoRR*, abs/2207.00747.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, and Denny Zhou. 2022e. Self-consistency improves chain of thought reasoning in language models. *CoRR*, abs/2203.11171.

Yan Wang, Xiaojiang Liu, and Shuming Shi. 2017. Deep neural solver for math word problems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 845–854, Copenhagen, Denmark. Association for Computational Linguistics.

Taylor W. Webb, Keith J. Holyoak, and Hongjing Lu. 2022. Emergent analogical reasoning in large language models. *CoRR*, abs/2212.09196.

Albert Webson and Ellie Pavlick. 2022. Do prompt-based models really understand the meaning of their prompts? In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 2300–2344. Association for Computational Linguistics.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022a. Emergent abilities of large language models. *Transactions on Machine Learning Research*. Survey Certification.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. 2022b. Chain of thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, December 6-14, 2022, virtual*.

Yixuan Weng, Minjun Zhu, Shizhu He, Kang Liu, and Jun Zhao. 2022. Large language models are reasoners with self-verification. *CoRR*, abs/2212.09561.

Sarah Wiegreffe, Jack Hessel, Swabha Swayamdipta, Mark O. Riedl, and Yejin Choi. 2022. Reframing human-ai collaboration for generating free-text explanations. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 632–658. Association for Computational Linguistics.

Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023. Visual chatgpt: Talking, drawing and editing with visual foundation models. *CoRR*, abs/2303.04671.

Tongshuang Wu, Ellen Jiang, Aaron Donsbach, Jeff Gray, Alejandra Molina, Michael Terry, and Carrie J. Cai. 2022a. Promptchainer: Chaining large language model prompts through visual programming. In *CHI '22: CHI Conference on Human Factors in Computing Systems, New Orleans, LA, USA, 29 April 2022 - 5 May 2022, Extended Abstracts*, pages 359:1–359:10. ACM.

Yuhuai Wu, Albert Q. Jiang, Wenda Li, Markus N. Rabe, Charles Staats, Mateja Jamnik, and Christian Szegedy. 2022b. Autoformalization with large language models. *CoRR*, abs/2205.12615.

Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2022. An explanation of in-context learning as implicit bayesian inference. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Frank F. Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. 2022. A systematic evaluation of large language models of code. In *MAPS@PLDI 2022: 6th ACM SIGPLAN International Symposium on Machine Programming, San Diego, CA, USA, 13 June 2022*, pages 1–10. ACM.

Jingjing Xu, Wangchunshu Zhou, Zhiyi Fu, Hao Zhou, and Lei Li. 2021. A survey on green deep learning. *CoRR*, abs/2111.05193.

Jingfeng Yang, Haoming Jiang, Qingyu Yin, Danqing Zhang, Bing Yin, and Diyi Yang. 2022a. SEQZERO: few-shot compositional semantic parsing with sequential prompts and zero-shot models. In *Findings of the Association for Computational Linguistics: NAACL 2022, Seattle, WA, United States, July 10-15,*

*2022*, pages 49–60. Association for Computational Linguistics.

Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023. MM-REACT: prompting chatgpt for multimodal reasoning and action. *CoRR*, abs/2303.11381.

Zhicheng Yang, Jinghui Qin, Jiaqi Chen, Liang Lin, and Xiaodan Liang. 2022b. LogicSolver: Towards interpretable math word problem solving with logical prompt-enhanced learning. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1–13, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei. 2022c. Language models as inductive reasoners. *CoRR*, abs/2212.10923.

Fu Yao, Peng Hao, and Khot Tushar. 2022. How does gpt obtain its ability? tracing emergent abilities of language models to their sources. *Yao Fu's Notion*.

Xi Ye and Greg Durrett. 2022. The unreliability of explanations in few-shot prompting for textual reasoning. In *Advances in Neural Information Processing Systems*.

Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. 2023. Large language models are versatile decomposers: Decompose evidence and questions for table-based reasoning. *CoRR*, abs/2301.13808.

Ori Yoran, Tomer Wolfson, Ben Bogin, Uri Katz, Daniel Deutch, and Jonathan Berant. 2023. Answering questions by meta-reasoning over multiple chains of thought.

Nathan Young, Qiming Bao, Joshua Bensemann, and Michael Witbrock. 2022. Abductionrules: Training transformers to explain unexpected inputs. In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 218–227. Association for Computational Linguistics.

Ping Yu, Tianlu Wang, Olga Golovneva, Badr AlKhamissy, Gargi Ghosh, Mona T. Diab, and Asli Celikyilmaz. 2022. ALERT: adapting language models to reasoning tasks. *CoRR*, abs/2212.08286.

Weizhe Yuan and Pengfei Liu. 2022. restructured pre-training. *CoRR*, abs/2206.11147.

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. STar: Bootstrapping reasoning with reasoning. In *Advances in Neural Information Processing Systems*.

Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. From recognition to cognition: Visual commonsense reasoning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 6720–6731. Computer Vision Foundation / IEEE.

Hanlin Zhang, YiFan Zhang, Li Erran Li, and Eric Xing. 2022. The impact of symbolic representations on in-context learning for few-shot reasoning. In *NeurIPS 2022 Workshop on Neuro Causal and Symbolic AI (nCSI)*.

Ningyu Zhang, Lei Li, Xiang Chen, Xiaozhuan Liang, Shumin Deng, and Huajun Chen. 2023a. Multimodal analogical reasoning over knowledge graphs. In *The Eleventh International Conference on Learning Representations*.

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2023b. Automatic chain of thought prompting in large language models. In *The Eleventh International Conference on Learning Representations*.

Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023c. Multimodal chain-of-thought reasoning in language models. *CoRR*, abs/2302.00923.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A survey of large language models. *CoRR*, abs/2303.18223.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations*.

Hattie Zhou, Azade Nova, Hugo Larochelle, Aaron C. Courville, Behnam Neyshabur, and Hanie Sedghi. 2022a. Teaching algorithmic reasoning via in-context learning. *CoRR*, abs/2211.09066.

Pei Zhou, Andrew Zhu, Jennifer Hu, Jay Pujara, Xiang Ren, Chris Callison-Burch, Yejin Choi, and Prithviraj Ammanabrolu. 2022b. An AI dungeon master's guide: Learning to converse and guide with intents and theory-of-mind in dungeons and dragons. *CoRR*, abs/2212.10060.

## A Appendix

### A.1 Related Survey

As this area is relatively nascent, only a few surveys exist. Closest to our work, Huang and Chang (2022) gives a survey towards reasoning with large language models. Dong et al. (2023) organizes and discusses the advanced techniques of in-context learning. Zhao et al. (2023) reviews the latest advancements in Large Language Models (LLMs)
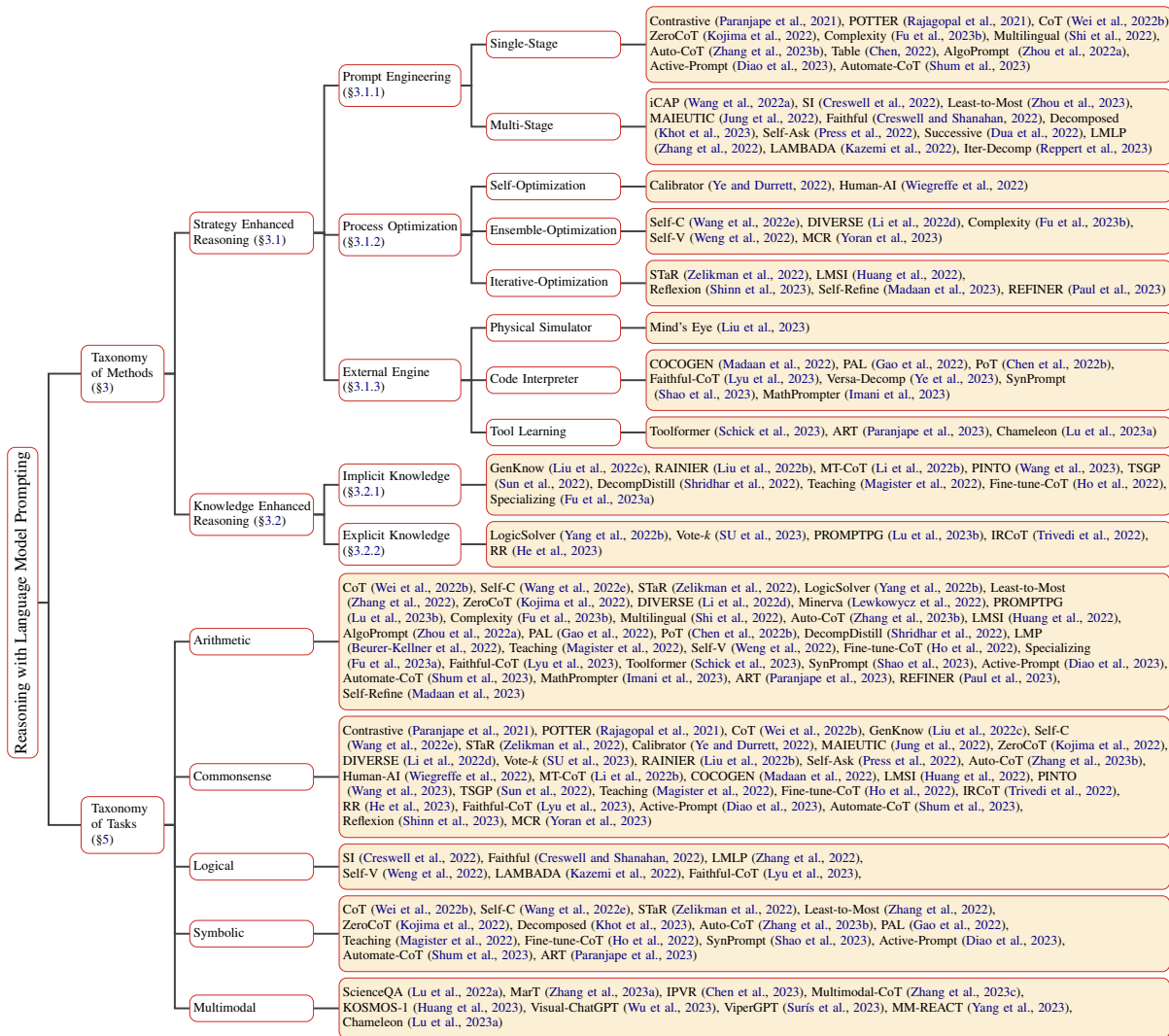
Figure 8: Taxonomy of Reasoning with Language Model Prompting.

and delves into the unresolved challenges that will shape future developments. Bhargava and Ng (2022) covers methods for commonsense knowledge reasoning and generation with pre-trained LMs. Lu et al. (2022b) reviews the key tasks, datasets, and methods at the intersection of mathematical reasoning and deep learning over the past decade. Liang et al. (2022a) surveys knowledge graph reasoning tracing from static to temporal and then to multi-modal knowledge graphs. Mialon et al. (2023) reviews works in which language models (LMs) are augmented with reasoning skills and the ability to use tools. Hamilton et al. (2022) conducts a survey of studies implementing neural-symbolic (NeSy) NLP approaches for reasoning and so on. Guo et al. (2022) provides a survey of several popular works dealing with uncertainty reasoning. Qin et al. (2023) concentrates on the leverage of external tools by LLMs which is also called Tool Learning. Other surveys focusing on prompt learning (Liu et al., 2022d) or pre-trained models (Qiu et al., 2020; Du et al., 2022) are also related to our work.

Unlike those surveys, in this paper, we conduct a review of reasoning with LM prompting, hoping to systematically understand the methodologies, compare different methods and inspire new ideas.

## A.2 Taxonomy of Methods and Tasks

We list the complete taxonomy of reasoning with language model prompting from methods and tasks in Figure 8.

## A.3 Performance Comparison of LMs with Different Scales

To show the generalization of discussions in §4.1 on different reasoning tasks, we additionally show the performance comparison of LMs with different scales on CommonsenseQA (Talmor et al., 2019) of commonsense reasoning in Figure 9.

## A.4 Detailed Information of Reasoning Benchmarks

In § 5, we give a brief overview on benchmarks and tasks requiring various reasoning skills. We list more benchmarks and show their key statistics in Table 2. Apart from the above-mentioned specific reasoning tasks in § 5, there are some benchmarks (Lake and Baroni, 2017; Srivastava et al., 2022; Yu et al., 2022) that can evaluate the model's more diverse and generalized reasoning capabilities, which


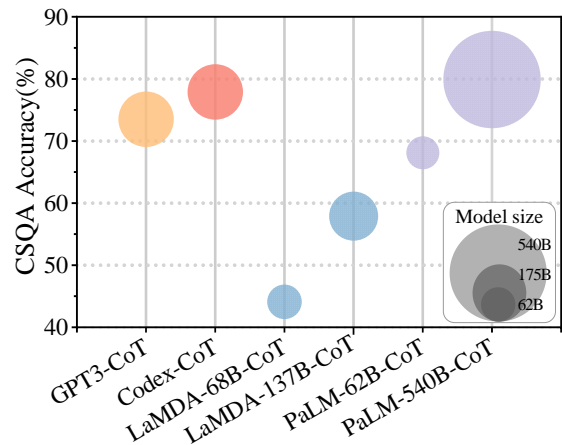
Figure 9: Performance of different language model scales on commonsense reasoning. Representatively, We show CoT (Wei et al., 2022b) experimental results on CommonsenseQA (Talmor et al., 2019).

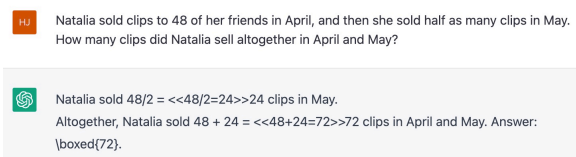can also be included in the category of reasoning tasks.

## A.5 Reasoning with ChatGPT



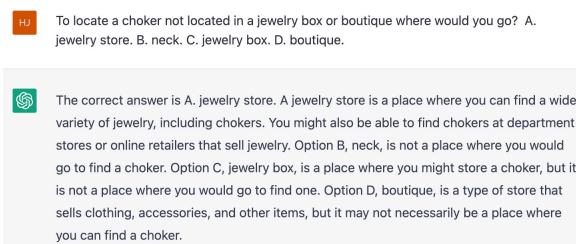Figure 10: A test case from GSM8K (Cobbe et al., 2021) on ChatGPT (OpenAI, 2022).



Figure 11: A test case from CommonsenseQA (Talmor et al., 2019) on ChatGPT (OpenAI, 2022).

Recently, OpenAI (2022) develops ChatGPT, an AI chatbot system that has attracted tremendous users. ChatGPT is trained on a massive dataset of text and is able to generate human-like responses to a wide variety of prompts, the promising approach for which is called Reinforcement Learning from Human Feedback (Ouyang et al., 2022). The backbone of ChatGPT is from a model in the GPT-3.5

| Task | Dataset | Size | | | |
|------|---------|------|------|------|------|
| | | **Train** | **Valid** | **Test** | **All** |
| Arithmetic Reasoning | AddSub (Hosseini et al., 2014) | 395 | - | - | 395 |
| | SingleOp (Roy et al., 2015) | 562 | - | - | 562 |
| | SingleEq (Koncel-Kedziorski et al., 2015) | 508 | - | - | 508 |
| | MultiArith (Roy and Roth, 2015) | 600 | - | - | 600 |
| | Dophin18k (Huang et al., 2016) | 18,460 | - | - | 18,460 |
| | MAWPS (Koncel-Kedziorski et al., 2016) | 1,921 | - | - | 1,921 |
| | Math23k (Wang et al., 2017) | 23,161 | - | - | 23,161 |
| | AQUA-RAT (Ling et al., 2017b) | 97,467 | - | 254 | 97,721 |
| | MathQA (Amini et al., 2019) | 29,807 | 4,471 | 2,981 | 37,259 |
| | DROP (Dua et al., 2019) | 5,850 | - | - | 5,850 |
| | ASDiv (Miao et al., 2020) | 1,217 | - | - | 1,217 |
| | GSM8K (Cobbe et al., 2021) | 7,473 | - | 1,319 | 8,792 |
| | SVAMP (Patel et al., 2021) | 1,000 | - | - | 1,000 |
| | MATH (Hendrycks et al., 2021) | 7,500 | - | 5,000 | 12,500 |
| | NumGLUE (Mishra et al., 2022b) | 101,835 | - | - | 101,835 |
| | Lila (Mishra et al., 2022a) | 133,815 | - | - | 133,815 |
| Symbolic Reasoning | Last Letter Concatenation (Wei et al., 2022b) | - | - | - | - |
| | Coin Flip (Wei et al., 2022b) | - | - | - | - |
| | Reverse List (Wei et al., 2022b) | - | - | - | - |
| Commonsense Reasoning | ARC (Clark et al., 2018) | 3,370 | 869 | 3,548 | 7,787 |
| | OpenBookQA (Mihaylov et al., 2018) | 4,957 | 500 | 500 | 5,957 |
| | CommonsenseQA (Talmor et al., 2019) | 9,741 | 1,221 | 1,140 | 12,102 |
| | PIQA (Bisk et al., 2020) | 16,000 | 2,000 | 3,000 | 21,000 |
| | StrategyQA (Geva et al., 2021) | 2,290 | - | 490 | 2,780 |
| Logical Reasoning | RuleTaker (Clark et al., 2020) | 14,135 | 2,019 | 3,038 | 20,192 |
| | ProofWriter (Tafjord et al., 2021) | - | - | - | - |
| | EntailmentBank (Dalvi et al., 2021) | 1,313 | 187 | 340 | 1,840 |
| | CLUTRR (Sinha et al., 2019) | 6,016 | - | - | 6,016 |
| | DEER (Yang et al., 2022c) | 1,200 | - | - | 1,200 |
| Multimodal Reasoning | VCR (Zellers et al., 2019) | 212,923 | 26,534 | 25,263 | 264,720 |
| | VisualCOMET (Park et al., 2020) | 1,174,063 | 146,332 | 145,309 | 1,465,704 |
| | VLEP (Lei et al., 2020) | 20,142 | 4,392 | 4,192 | 28,726 |
| | PMR (Dong et al., 2022) | 12,080 | 1,538 | 1,742 | 15,360 |
| | ScienseQA (Lu et al., 2022a) | 12,726 | 4,241 | 4,241 | 21,208 |

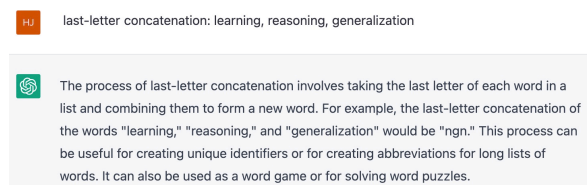Table 2: An overview of benchmarks and tasks on reasoning.



Figure 12: A test case from Last Letter Concatenation (Wei et al., 2022b) on ChatGPT (OpenAI, 2022).

large LM series[8]. In order to savor the reasoning ability of large LMs more realistically, we conduct some case tests on ChatGPT. Concretely, we pick out a piece of data from GSM8K (Cobbe et al., 2021), CommonsenseQA (Talmor et al., 2019) and Last Letter Concatenation (Wei et al., 2022b) which respectively represent arithmetic reasoning, commonsense reasoning, and symbolic reasoning.

Then we test each of the selected data on ChatGPT directly. Results can be seen in Figure 10-12.
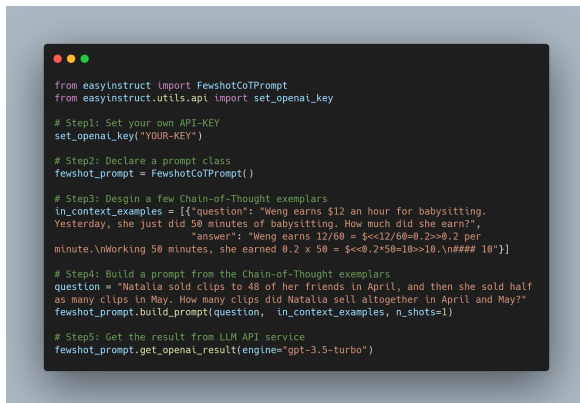
Figure 10 shows that given a math problem in GSM8K (Cobbe et al., 2021), ChatGPT outputs a reasoning process and a correct answer without in-context exemplars. This blazes its powerful arithmetic reasoning ability. The reasoning process has the same format as the gold label in GSM8K, indicating that GSM8K may be contained in the training corpus of ChatGPT.

In Figure 11, we test ChatGPT on a piece of data in CommsonsenseQA (Talmor et al., 2019). It not only gives the correct answer but additionally details why each option is right or wrong, which does not appear in the gold label of the dataset. This demonstrates the strong commonsense reasoning ability of ChatGPT.

Figure 12 is a case in Last Letter Concatenation (Wei et al., 2022b). We observe that although Chat-

GPT gives a detailed and accurate description of last letter concatenation, it fails to answer the given question, showing that its symbolic reasoning capability is not as excellent as the above two.

## A.6 Reasoning using EasyInstruct



```python
from easyinstruct import FewshotCoTPrompt
from easyinstruct.utils.api import set_openai_key

# Step1: Set your own API-KEY
set_openai_key("YOUR-KEY")

# Step2: Declare a prompt class
fewshot_prompt = FewshotCoTPrompt()

# Step3: Desgin a few Chain-of-Thought exemplars
in_context_examples = [{"question": "Weng earns $12 an hour for babysitting.
Yesterday, she just did 50 minutes of babysitting. How much did she earn?",
                        "answer": "Weng earns 12/60 = $<<12/60=0.2>>0.2 per
minute.\nWorking 50 minutes, she earned 0.2 x 50 = $<<0.2*50=10>>10.\n#### 10"}]

# Step4: Build a prompt from the Chain-of-Thought exemplars
question = "Natalia sold clips to 48 of her friends in April, and then she sold half
as many clips in May. How many clips did Natalia sell altogether in April and May?"
fewshot_prompt.build_prompt(question,  in_context_examples, n_shots=1)

# Step5: Get the result from LLM API service
fewshot_prompt.get_openai_result(engine="gpt-3.5-turbo")
```

Figure 13: A test case from GSM8K (Cobbe et al., 2021) using EasyInstruct (Ou et al., 2023).

## ACL 2023 Responsible NLP Checklist

### A   For every submission:

☑ A1. Did you describe the limitations of your work?
*Section 8 (Limitations)*

☐ A2. Did you discuss any potential risks of your work?
*Not applicable. Left blank.*

☑ A3. Do the abstract and introduction summarize the paper's main claims?
*abstract: Abstract; Introduction: section 1*

☑ A4. Have you used AI writing assistants when working on this paper?
*https://www.grammarly.com/*

### B   ☒ Did you use or create scientific artifacts?

*Left blank.*

☐ B1. Did you cite the creators of artifacts you used?
*Not applicable. Left blank.*

☐ B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
*Not applicable. Left blank.*

☐ B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
*Not applicable. Left blank.*

☐ B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
*Not applicable. Left blank.*

☐ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
*Not applicable. Left blank.*

☐ B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
*Not applicable. Left blank.*

### C   ☒ Did you run computational experiments?

*Left blank.*

☐ C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
*Not applicable. Left blank.*

☐ C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?
*Not applicable. Left blank.*

☐ C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?
*Not applicable. Left blank.*

☐ C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?
*Not applicable. Left blank.*

**D  ☒ Did you use human annotators (e.g., crowdworkers) or research with human participants?**

*Left blank.*

☐ D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?
*Not applicable. Left blank.*

☐ D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?
*Not applicable. Left blank.*

☐ D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?
*Not applicable. Left blank.*

☐ D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?
*Not applicable. Left blank.*

☐ D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?
*Not applicable. Left blank.*