

Table 4: Standard prompting versus chain of thought prompting on five commonsense reasoning benchmarks. Chain of thought prompting is an emergent ability of model scale—it does not positively impact performance until used with a model of sufficient scale.

Model		CSQA		StrategyQA		Date		Sports		SayCan	
		standard	CoT	standard	CoT	standard	CoT	standard	CoT	standard	CoT
UL2	20B	34.2	51.4	59.0	53.3	13.5	14.0	57.9	65.3	20.0	41.7
LaMDA	420M	20.1	19.2	46.4	24.9	1.9	1.6	50.0	49.7	7.5	7.5
	2B	20.2	19.6	52.6	45.2	8.0	6.8	49.3	57.5	8.3	8.3
	8B	19.0	20.3	54.1	46.8	9.5	5.4	50.0	52.1	28.3	33.3
	68B	37.0	44.1	59.6	62.2	15.5	18.6	55.2	77.5	35.0	42.5
	137B	53.6	57.9	62.4	65.4	21.5	26.8	59.5	85.8	43.3	46.6
GPT	350M	14.7	15.2	20.6	0.9	4.3	0.9	33.8	41.6	12.5	0.8
	1.3B	12.0	19.2	45.8	35.7	4.0	1.4	0.0	26.9	20.8	9.2
	6.7B	19.0	24.0	53.6	50.0	8.9	4.9	0.0	4.4	17.5	35.0
	175B	79.5	73.5	65.9	65.4	43.8	52.1	69.6	82.4	81.7	87.5
Codex	-	82.3	77.9	67.1	73.2	49.0	64.8	71.7	98.5	85.8	88.3
PaLM	8B	19.8	24.9	55.6	53.5	12.9	13.1	55.1	75.2	34.2	40.0
	62B	65.4	68.1	58.4	63.4	29.8	44.7	72.1	93.6	65.8	70.0
	540B	78.1	79.9	68.6	77.8	49.0	65.3	80.5	95.4	80.8	91.7

Table 5: Standard prompting versus chain of thought prompting enables length generalization to longer inference examples on two symbolic manipulation tasks.

Model		Last Letter Concatenation				Coin Flip (state tracking)					
		2		OOD: 3		OOD: 4		2		OOD: 3	
		standard	CoT	standard	CoT	standard	CoT	standard	CoT	standard	CoT
UL2	20B	0.6	18.8	0.0	0.2	0.0	0.0	70.4	67.1	51.6	52.2
LaMDA	420M	0.3	1.6	0.0	0.0	0.0	0.0	52.9	49.6	50.0	50.5
	2B	2.3	6.0	0.0	0.0	0.0	0.0	54.9	55.3	47.4	48.7
	8B	1.5	11.5	0.0	0.0	0.0	0.0	52.9	55.5	48.2	49.6
	68B	4.4	52.0	0.0	0.8	0.0	2.5	56.2	83.2	50.4	69.1
	137B	5.8	77.5	0.0	34.4	0.0	13.5	49.0	99.6	50.7	91.0
PaLM	8B	2.6	18.8	0.0	0.0	0.0	0.2	60.0	74.4	47.3	57.1
	62B	6.8	85.0	0.0	59.6	0.0	13.4	91.4	96.8	43.9	91.0
	540B	7.6	99.4	0.2	94.8	0.0	63.0	98.1	100.0	49.3	98.6

Table 6: Ablation and robustness results for arithmetic reasoning datasets. Chain of thought generally outperforms ablations by a large amount. “Equation only” performs in between standard prompting and chain of thought prompting, as it allows for intermediate reasoning steps via equations but does not leverage natural language. Chain of thought prompting has variance (as expected) when used with prompts written by different annotators or when using other exemplars, but still outperforms standard prompting by a large margin. Standard deviation shown is for different order of few-shot prompting exemplars, with five different random seeds. Results here are shown for LaMDA 137B, as additional queries for GPT-3 and PaLM are both limited and expensive.

	GSM8K	SVAMP	ASDiv	MAWPS
Standard prompting	6.5 ± 0.4	29.5 ± 0.6	40.1 ± 0.6	43.2 ± 0.9
Chain of thought prompting	14.3 ± 0.4	36.7 ± 0.4	46.6 ± 0.7	57.9 ± 1.5
Ablations				
· equation only	5.4 ± 0.2	35.1 ± 0.4	45.9 ± 0.6	50.1 ± 1.0
· variable compute only	6.4 ± 0.3	28.0 ± 0.6	39.4 ± 0.4	41.3 ± 1.1
· reasoning after answer	6.1 ± 0.4	30.7 ± 0.9	38.6 ± 0.6	43.6 ± 1.0
Robustness				
· different annotator (B)	15.5 ± 0.6	35.2 ± 0.4	46.5 ± 0.4	58.2 ± 1.0
· different annotator (C)	17.6 ± 1.0	37.5 ± 2.0	48.7 ± 0.7	60.1 ± 2.0
· intentionally concise style	11.1 ± 0.3	38.7 ± 0.8	48.0 ± 0.3	59.6 ± 0.7
· exemplars from GSM8K (α)	12.6 ± 0.6	32.8 ± 1.1	44.1 ± 0.9	53.9 ± 1.1
· exemplars from GSM8K (β)	12.7 ± 0.5	34.8 ± 1.1	46.9 ± 0.6	60.9 ± 0.8
· exemplars from GSM8K (γ)	12.6 ± 0.7	35.6 ± 0.5	44.4 ± 2.6	54.2 ± 4.7

Table 7: Ablation and robustness results for four datasets in commonsense and symbolic reasoning. Chain of thought generally outperforms ablations by a large amount. Chain of thought prompting has variance (as expected) when used with prompts written by different annotators or when using other exemplars, but still outperforms standard prompting by a large margin. Standard deviation shown is for different order of few-shot prompting exemplars, with five different random seeds. Results here are shown for LaMDA 137B, as additional queries for GPT-3 and PaLM are both limited and expensive. The exception is that we run SayCan using PaLM here, as the SayCan evaluation set is only 120 examples and therefore less expensive to run multiple times.

	Commonsense			Symbolic	
	Date	Sports	SayCan	Concat	Coin
Standard prompting	21.5 ± 0.6	59.5 ± 3.0	80.8 ± 1.8	5.8 ± 0.6	49.0 ± 2.1
Chain of thought prompting	26.8 ± 2.1	85.8 ± 1.8	91.7 ± 1.4	77.5 ± 3.8	99.6 ± 0.3
Ablations					
· variable compute only	21.3 ± 0.7	61.6 ± 2.2	74.2 ± 2.3	7.2 ± 1.6	50.7 ± 0.7
· reasoning after answer	20.9 ± 1.0	63.0 ± 2.0	83.3 ± 0.6	0.0 ± 0.0	50.2 ± 0.5
Robustness					
· different annotator (B)	27.4 ± 1.7	75.4 ± 2.7	88.3 ± 1.4	76.0 ± 1.9	77.5 ± 7.9
· different annotator (C)	25.5 ± 2.5	81.1 ± 3.6	85.0 ± 1.8	68.1 ± 2.2	71.4 ± 11.1

C Extended Related Work

Chain-of-thought prompting is a general approach that is inspired by several prior directions: prompting, natural language explanations, program synthesis/execution, numeric and logical reasoning, and intermediate language steps.

C.1 Prompting

The recent success of large-scale language models has led to growing interest in improving their capability to perform tasks via prompting (Brown et al. (2020), and see Liu et al. (2021) for a survey). This paper falls in the category of general prompting approaches, whereby input prompts are optimized to allow a single large language model to better perform a variety of tasks (Li and Liang, 2021; Lester et al., 2021; Reif et al., 2022, *inter alia*).

One recent line of work aims to improve the ability of language models to perform a task by providing instructions that describe the task (Raffel et al., 2020; Wei et al., 2022a; Ouyang et al., 2022; Sanh et al., 2022; Wang et al., 2022b). This line of work is related because it also augments input–output pairs with meta-data. But whereas an instruction augments the input to a task (instructions are typically prepended to the inputs), chain-of-thought prompting augments the outputs of language models. Another related direction is sequentially combining the outputs of language models; human–computer interaction (HCI) work (Wu et al., 2022a,b) has shown that combining sequential generations of language models improves task outcomes in a 20-person user study.

C.2 Natural language explanations

Another closely related direction uses natural language explanations (NLEs), often with the goal of improving model interpretability (Zhou et al., 2020; Wiegreffe and Marasović, 2021, *inter alia*). That line of work typically focuses on natural language inference (Camburu et al., 2018; Yordanov et al., 2021; Bostrom et al., 2021), and produces explanations either simultaneously to or after the final prediction (Narang et al., 2020; Majumder et al., 2021; Wiegreffe et al., 2021, 2022). By contrast, the chain of thought processing considered in this paper occurs *before* the final answer. And while NLE aims mostly to improve neural network interpretability (Rajagopal et al., 2021), the goal of chain-of-thought prompting is to allow models to decompose multi-hop reasoning tasks into multiple steps—interpretability is just a side effect. Marasović et al. (2022) show that prompt-based finetuning with NLE improves NLI and classification performance, though they largely focus on evaluating explanation plausibility. In comparison, our work focuses on a range of arithmetic, commonsense, and symbolic tasks that require multi-hop reasoning.

C.3 Program synthesis and execution

Using intermediate reasoning steps has a long history in program synthesis and execution (Zaremba and Sutskever, 2014, *inter alia*). Recent work along in this direction has included a number of architectural innovations (Cai et al., 2017; Dong et al., 2019; Yan et al., 2020), as well as the use of large language models (Chen et al., 2021; Austin et al., 2021). The program execution work closest to ours is perhaps Nye et al. (2021), which show that large language models can perform up to 10-digit addition, evaluate polynomials, and execute python programs. Whereas generating a program and then executing it can be viewed as a type of reasoning, our work generalizes such domain-specific primitives to natural language, which is open-domain and relevant to any text-to-text NLP task in principle.

C.4 Numeric and logical reasoning

Numeric and logical reasoning has been a long-studied task in machine learning and natural language processing (Lev et al., 2004, *inter alia*). Recent work has also aimed to inject numeric reasoning abilities in language models in various ways, such as augmenting BERT with a predefined set of executable operations (Andor et al., 2019), including a graph neural network (Ran et al., 2019), and using specialized training procedures (Piękos et al., 2021). Another line of work aims to enable language models to perform logical or formal reasoning, often by verbalizing the rules in natural language formal rules using language (Clark et al., 2020; Saeed et al., 2021; Liang et al., 2021).