# Not All Votes Count! Programs as Verifiers Improve Self-Consistency of Language Models for Math Reasoning

**Vernon Y.H. Toh[1], Deepanway Ghosal[1*], Soujanya Poria[1]**

[1] Singapore University of Technology and Design

## Abstract

Large language models (LLMs) have shown increasing proficiency in solving mathematical reasoning problems. However, many current open-source LLMs often still make calculation and semantic understanding errors in their intermediate reasoning steps. In this work, we propose PROVE, a simple yet effective framework that uses program-based verification as a heuristic to filter out potentially incorrect reasoning paths before aggregating the final answers. Instead of relying on vanilla majority voting, our approach rejects solutions whose corresponding program outputs are inconsistent with the generated solution, aggregating only those validated by Python programs. We conducted extensive experiments on 13 open-source LLMs from various model families and sizes, ranging from 0.5B to 13B parameters, across seven math benchmarks. We demonstrate that PROVE consistently outperforms vanilla majority voting as a heuristic for solving mathematical reasoning tasks across all datasets and model sizes. Notably, PROVE increases accuracy on the GSM8K benchmark from 48.85% to 53.83% for Qwen2-0.5B-Instruct, from 65.66% to 73.01% for Llama-3.2-1B-Instruct, from 73.39% to 79.61% for Gemma-2-2b-it, and from 41.32% to 59.51% for Llama-2-7B-chat. Our codes are available at https://github.com/declare-lab/prove.

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable success in mathematical reasoning tasks, especially with advanced models like GPT-4 (OpenAI et al., 2024). However, smaller open-source LLMs, such as Mistral, Gemma 2, and Llama 2 (Jiang et al., 2023; Team et al., 2024; Touvron et al., 2023) fall short in mathematical reasoning, often producing inaccurate content and
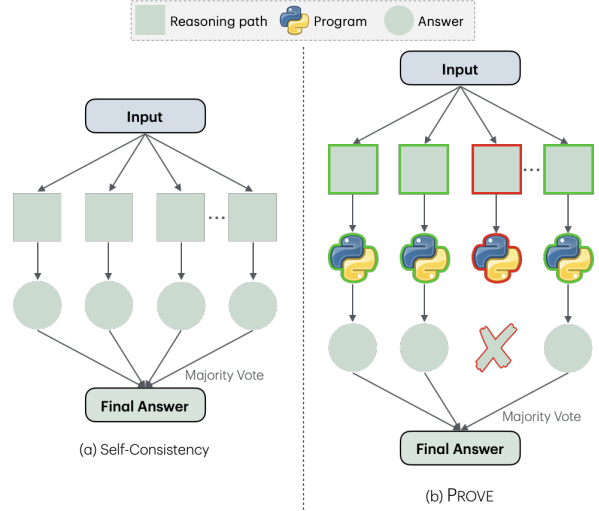
---
*Now at Deepmind.



Figure 1: Comparison of self-consistency and PROVE.

failing to handle complex calculations. Many work since then have aimed to tackle those challenges by using different prompting methods such as Plan-and-Solve (PS) Prompting (Wang et al., 2023a) which aims to mitigate calculation errors and improve the overall quality of the generated reasoning by adopting a two-step process: first, devise a plan to divide the entire task into smaller subtasks and then carrying out the subtasks according to the plan. Non-prompting approaches, such as fine-tuning LLMs for reasoning tasks, have also shown performance improvements. However this typically requires large amounts of distilled data from models like GPT-4 (Wang et al., 2024; Gou et al., 2024), which comes with high computational and generation costs.

Despite these advancements, fundamental challenges remain (Banerjee et al., 2024). Autoregressive models lack mechanisms to correct their own errors, meaning that once a mistake is made, it can propagate through the entire response. As a result, LLMs are still vulnerable to hallucinations which can lead to incorrect reasoning or calcula-

**#1:** Qwen2-0.5B-Instruct    **#2:** Qwen2-1.5B-Instruct    **#3:** Qwen2-7B-Instruct    **#4:** Gemma-2-2B-it    **#5:** Gemma-2-9B-it    **#6:** Phi-3-mini-4k-instruct    **#7:** Mistral-7B-Instruct-v0.3
**#8:** Llama-2-7B-chat    **#9:** Llama-2-13B-chat    **#10:** Llama-3-8B-Instruct    **#11:** Llama-3.1-8B-Instruct    **#12:** Llama-3.2-1B-Instruct    **#13:** Llama-3.2-3B-Instruct
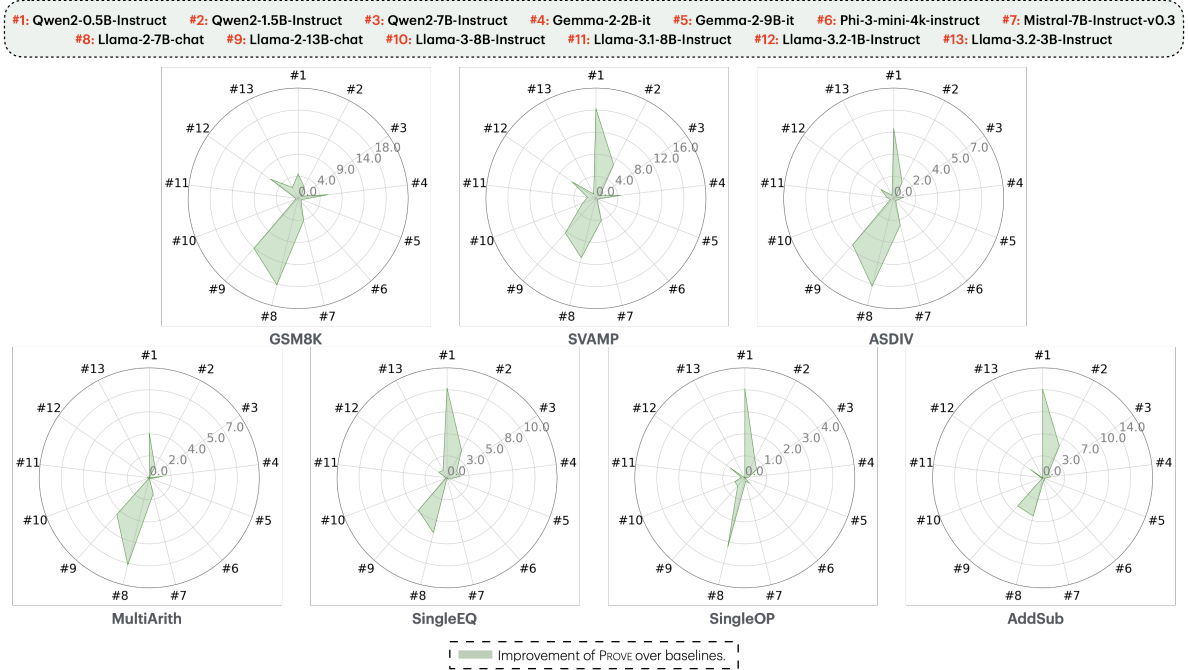
Figure 2: The improvement with PROVE when used with different LLMs on seven benchmarks.

tion mistakes. Given these limitations, a promising direction is to enhance reasoning by leveraging the knowledge already embedded within LLMs themselves (Wang et al., 2023b; Hao et al., 2023; Madaan et al., 2023; Weng et al., 2023; Wang et al., 2023a). Following Wang et al. (2023b), one can sample multiple reasoning paths to solve a question and select the most common answer. While aggregating answers like this improves performance, sometimes the most common answer could still be wrong. Recent efforts have tried to develop better heuristics for selecting answers from multiple reasoning paths, such as verification via formalization (Lin et al., 2024). However, this approach requires few-shot demonstration examples and more capable models, like GPT-3.5, to effectively verify solutions through autoformalization.

In this paper, we introduce PROVE, a *frustratingly simple* approach that utilizes **PRO**grams as **VE**rifiers. Our framework is model-agnostic and eliminates the need for LLM fine-tuning or few-shot prompting. Instead, it leverages the capabilities of current open-source LLMs to translate natural language plans and solutions into Python programs for verification. Unlike methods that rely on simple majority voting across all generated natural language solutions, we aggregate only those solutions that are successfully verified by their corresponding Python programs. While a natural language plan and solution provides an abstract outline

of the key steps in solving a math word problem, verifying its correctness can be challenging. Python program implementations offer a more precise representation, making it easier to accurately identify correct natural language solutions.

We evaluate PROVE using 13 LLMs from various model families and sizes (Abdin et al., 2024; Touvron et al., 2023; Dubey et al., 2024; Jiang et al., 2023; Team et al., 2024; Yang et al., 2024), ranging from 0.5B to 13B parameters across 7 mathematical reasoning datasets (Cobbe et al., 2021; Patel et al., 2021; Miao et al., 2020; Roy and Roth, 2015; Koncel-Kedziorski et al., 2015; Roy et al., 2015; Hosseini et al., 2014). The results demonstrate that PROVE consistently outperforms vanilla majority voting (Wang et al., 2023b) across all model sizes, evaluated on various math reasoning.

**Contributions.** We summarize the contributions of our paper as follows:

- We propose a novel framework PROVE, that utilizes programs as verifiers, serving as a heuristic to filter out potentially incorrect reasoning paths before aggregating the answers.
- Our experimental results demonstrate that PROVE significantly outperforms baseline methods across seven mathematical reasoning datasets.
- We present comprehensive ablation studies and qualitative analyses that clearly demonstrate the effectiveness of our framework.
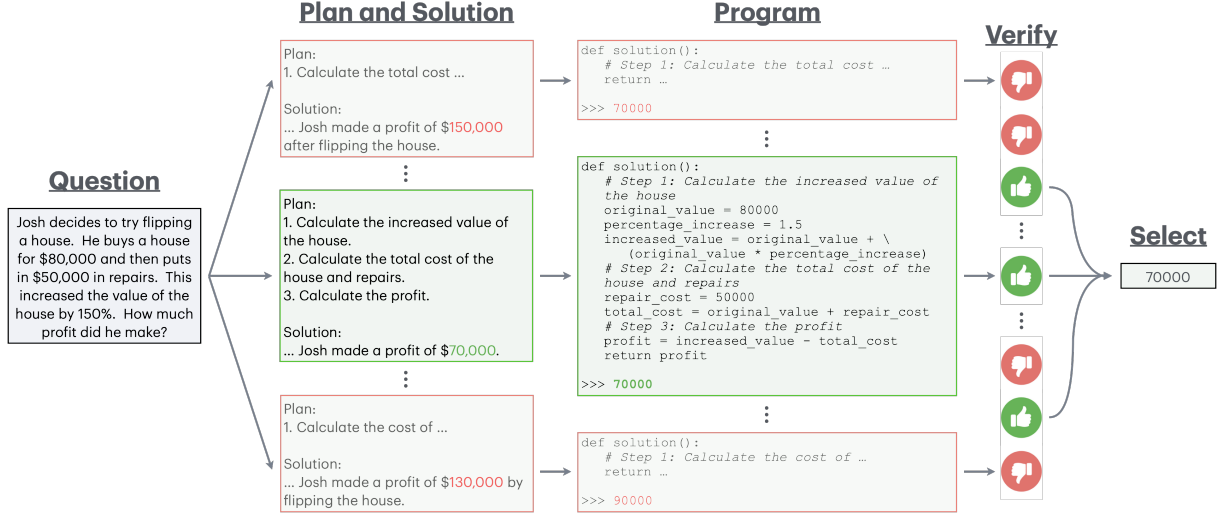
Figure 3: The PROVE Framework.

## 2 Method

**Overview.** As illustrated in Figure 3, PROVE framework starts by prompting an LLM to generate a plan and solution to solve a given math word problem. Next, we prompt an LLM to translate the generated plan and solution to a Python program. We then execute the Python program and compare its output to the generated solution. If the program's output matches the solution, it is considered valid, otherwise, the solution is filtered out. Finally, the remaining valid solutions are undergo a majority voting process to determine the final answer. It's important to note that all prompts in the PROVE framework are zero-shot.

**Generating plan and solution.** The first step in our framework involves prompting an LLM to generate a natural language plan and solution to solve a math word problems. For each math word problem $x$, we prompt an LLM to generate a plan and solution $s$ using the zero-shot prompt proposed by Wang et al. (2023a): "*Let's first understand the problem and devise a plan to solve the problem. Then, let's carry out the plan and solve the problem step by step*". We sample multiple plans and solutions from the LLM, using a temperature of 0.7, to produce a set of candidate plans and solutions $\{s_1, ..., s_n\}$.

**Translating plan and solution to Python programs.** After obtaining the set of candidate plans and solutions $\{s_1, ..., s_n\}$, we use each plan and solution $s_i$ individually as input to an LLM, prompting it to generate a Python program $p_i$ that implements the plan to solve the given math word prob-

lem. This process results in a corresponding Python program for each candidate plan and solution, yielding the pairs $\{(s_1, p_1), ..., (s_n, p_n)\}$, which are then used for verification.

**Verification and selection.** Lastly, we execute the Python program $p_i$ to obtain the output $\hat{p}_i$, which we use to verify the answer $\hat{s}_i$ from the generated solution $s_i$. The answer $\hat{s}_i$ is extracted using the answer extraction prompting method proposed by Kojima et al. (2022): "*Therefore, the answer (arabic numerals) is*". We consider the answer $\hat{s}_i$ valid if it matches the output $\hat{p}_i$ from the Python program. Once all candidate plans and solutions are verified, we perform majority voting on the remaining valid answers to determine the final answer. If no valid answer remains, we perform majority voting over all invalid answers.

## 3 Experiments

### 3.1 Setup

**Models and Datasets.** PROVE is a model-agnostic framework designed to integrate with any LLM for mathematical reasoning tasks. We evaluated 13 LLMs from various model families and sizes, ranging from 0.5B to 13B parameters. Specifically, evaluated on Phi-3-mini-4k-instruct (3.8B) (Abdin et al., 2024), Llama-2-7B-chat (Touvron et al., 2023), Llama-2-13B-chat, Llama-3-8B-Instruct (Dubey et al., 2024), Llama-3.1-8B-Instruct, Llama-3.2-1B-Instruct, Llama-3.2-3B-Instruct, Mistral-7B-Instruct-v0.3 (Jiang et al., 2023), Gemma-2-2B-it (Team et al., 2024), Gemma-2-9B-it, Qwen2-0.5B-Instruct (Yang et al.,

---

**Prompt**

Question: Janet's ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at the farmers' market daily for $2 per fresh duck egg. How much in dollars does she make every day at the farmers' market?
Answer: *Let's first understand the problem and devise a plan to solve the problem. Then, let's carry out the plan to solve the problem step by step.*

---

**Output by Phi-3-mini-4k-instruct**

Plan:
1. Calculate the number of eggs Janet uses for breakfast and baking.
2. Subtract the eggs used from the total number of eggs laid by the ducks to find out the number of eggs sold at the farmers' market.
3. Multiply the number of eggs sold by the price per egg to find the total amount Janet makes at the farmers' market.

Step-by-step Solution:
1. Eggs used for breakfast and baking:
Janet eats 3 eggs for breakfast and bakes 4 muffins with eggs. So, she uses a total of 3 + 4 = 7 eggs per day for these two purposes.
2. Eggs sold at the farmers' market:
Janet's ducks lay 16 eggs per day. After using 7 eggs for breakfast and baking, she has 16 - 7 = 9 eggs left to sell at the farmers' market.
3. Calculate the total amount Janet makes at the farmers' market:
Janet sells each egg for $2. Therefore, she makes a total of 9 eggs * $2/egg = $18 per day at the farmers' market.

**(a) Generating Plan and Solution**

---

**Prompt**

Question: Janet's ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at the farmers' market daily for $2 per fresh duck egg. How much in dollars does she make every day at the farmers' market?
Answer: Let's first understand the problem and devise a plan to solve the problem. Then, let's carry out the plan to solve the problem step by step.

Plan:
1. Calculate the number of eggs Janet uses for breakfast and baking.
2. Subtract the eggs used from the total number of eggs laid by the ducks to find out the number of eggs sold at the farmers' market.
3. Multiply the number of eggs sold by the price per egg to find the total amount Janet makes at the farmers' market.

Step-by-step Solution:
1. Eggs used for breakfast and baking:
Janet eats 3 eggs for breakfast and bakes 4 muffins with eggs. So, she uses a total of 3 + 4 = 7 eggs per day for these two purposes.
2. Eggs sold at the farmers' market:
Janet's ducks lay 16 eggs per day. After using 7 eggs for breakfast and baking, she has 16 - 7 = 9 eggs left to sell at the farmers' market.
3. Calculate the total amount Janet makes at the farmers' market:
Janet sells each egg for $2. Therefore, she makes a total of 9 eggs * $2/egg = $18 per day at the farmers' market.

*# Let's carry out the plan and answer this question by implementing a solution() function.*
*def solution():*
    *# Let's write a Python program step by step, and then return the numeric answer*

---

**Output by Phi-3-mini-4k-instruct**

```
def solution():
    # Step 1: Calculate the number of eggs Janet uses for breakfast and baking
    eggs_for_breakfast = 3
    eggs_for_baking = 4
    total_eggs_used = eggs_for_breakfast + eggs_for_baking

    # Step 2: Subtract the eggs used from the total number of eggs laid by the ducks
    eggs_laid_per_day = 16
    eggs_sold_at_market = eggs_laid_per_day - total_eggs_used

    # Step 3: Calculate the total amount Janet makes at the farmers' market
    price_per_egg = 2
    total_amount_made = eggs_sold_at_market * price_per_egg

    # Return the total amount made at the farmers' market
    return total_amount_made
```

**(b) Implementing Python Program from Plan and Solution**

Figure 4: An actual example of the prompts used in our framework, along with the outputs generated by Phi3-mini-4k-instruct, for the GSM8K task.

2024), Qwen2-1.5B-Instruct, and Qwen2-7B-Instruct. We conducted evaluations across seven mathematical reasoning datasets, including GSM8K (Cobbe et al., 2021), SVAMP (Patel et al., 2021), ASDIV (Miao et al., 2020), MultiArith (Roy and Roth, 2015), SingleEQ (Koncel-Kedziorski et al., 2015), SingleOP (Roy et al., 2015), and AddSub (Hosseini et al., 2014).

**Implementation Details.** In PROVE, we sample a total of 16 plans and solutions using a temperature of 0.7. To extract the answer from the generated solutions, we rely on the Phi-3-mini-4k-instruct model. We also used Phi-3-mini-4k-instruct to translate the plan and solution into Python programs for verification. With only 3.8 billion parameters, Phi-3-mini-4k-instruct offers more efficient inference. Notably, when Phi-3-mini-4k-instruct

is being evaluated, it performs self-verification by generating not only the plan and solution but also translating the plan and solution to a Python program to verify the results. During the selection phase, we perform majority voting on the remaining valid answer. If no valid answer remains, we perform majority voting over all invalid answers.

**Baselines.** We compare PROVE with two decoding strategies: single sampling (greedy decoding) and multiple sampling (self-consistency decoding) (Wang et al., 2023b). For the comparison, we use three different prompting techniques: zero-shot CoT (Kojima et al., 2022), zero-shot PoT (Chen et al., 2023), and zero-shot PS (Wang et al., 2023a). Performance is evaluated based on accuracy, where a problem is considered correct only if the final answer matches the ground truth.

| Method | Qwen-2 | | | Gemma-2 | | Phi-3 | Mistral | Llama-2 | | Llama-3 | Llama-3.1 | Llama-3.2 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0.5B | 1.5B | 7B | 2B | 9B | 3.8B | 7B | 7B | 13B | 8B | 8B | 1B | 3B |
| **GSM8K** | | | | | | | | | | | | | |
| Direct | 37.00 | 61.56 | 88.02 | 67.32 | 88.63 | 85.67 | 59.06 | 30.55 | 42.84 | 79.38 | 85.67 | 48.75 | 79.23 |
| Zero-shot CoT | 32.22 | 57.71 | 88.25 | 67.63 | 89.69 | 86.35 | 58.83 | 29.80 | 41.24 | 80.74 | 87.49 | 50.72 | 81.88 |
| Zero-shot PoT | 5.23 | 2.81 | 63.61 | 50.80 | 80.89 | 84.38 | 48.52 | 8.49 | 26.99 | 73.46 | 84.38 | 42.15 | 33.06 |
| Zero-shot PS | 25.63 | 51.02 | 85.90 | 63.00 | 89.01 | 86.50 | 58.07 | 28.35 | 39.80 | 80.36 | 86.66 | 49.05 | 81.12 |
| Zero-shot CoT @maj16 | 48.85 | 71.57 | 91.66 | 72.94 | 91.05 | 91.96 | 75.66 | 38.97 | 54.44 | 87.72 | 92.49 | 66.11 | 87.87 |
| Zero-shot PoT @maj16 | 28.43 | 37.00 | 85.67 | 59.21 | 86.43 | 91.36 | 70.96 | 27.22 | 40.86 | 84.07 | 90.37 | 58.45 | 83.32 |
| Zero-shot PS @maj16 | 42.00 | 68.46 | 90.98 | 73.39 | 91.96 | 92.19 | 78.54 | 41.32 | 54.21 | 89.61 | 92.57 | 65.66 | 88.86 |
| PROVE | **53.83** | **74.22** | **92.42** | **79.61** | **92.72** | **93.10** | **83.24** | **59.51** | **68.08** | **90.14** | **93.19** | **73.01** | **91.36** |
| **SVAMP** | | | | | | | | | | | | | |
| Direct | 45.70 | 65.90 | 92.40 | 72.90 | 89.70 | 91.70 | 69.30 | 61.40 | 66.20 | 83.00 | 89.10 | 70.90 | 85.70 |
| Zero-shot CoT | 41.00 | 62.10 | 89.90 | 74.40 | 88.70 | 89.10 | 71.80 | 56.40 | 62.80 | 85.00 | 86.80 | 68.40 | 86.50 |
| Zero-shot PoT | 13.60 | 5.40 | 47.00 | 73.40 | 91.80 | 90.70 | 52.80 | 25.50 | 44.10 | 86.70 | 88.00 | 62.60 | 66.30 |
| Zero-shot PS | 35.60 | 61.00 | 90.70 | 73.60 | 89.80 | 89.80 | 73.20 | 56.00 | 57.30 | 83.40 | 88.30 | 65.60 | 87.50 |
| Zero-shot CoT @maj16 | 57.10 | 78.30 | 93.60 | 80.00 | 91.30 | 94.00 | 84.40 | 69.70 | 77.80 | 90.50 | 92.30 | 82.30 | 91.90 |
| Zero-shot PoT @maj16 | 47.80 | 39.20 | 74.50 | 79.30 | 92.50 | 94.70 | 80.50 | 59.70 | 66.00 | 90.60 | 91.50 | 76.80 | 92.00 |
| Zero-shot PS @maj16 | 56.60 | 79.60 | 93.70 | 81.70 | 91.40 | 93.60 | 84.70 | 71.30 | 77.80 | 91.20 | 92.90 | 81.50 | 92.80 |
| PROVE | **73.80** | **86.70** | **94.60** | **86.40** | **92.90** | **95.10** | **88.90** | **82.70** | **86.40** | **93.90** | **94.50** | **87.70** | **93.70** |
| **ASDIV** | | | | | | | | | | | | | |
| Direct | 60.02 | 75.43 | 91.94 | 81.30 | 91.70 | 93.32 | 73.57 | 61.31 | 67.51 | 85.26 | 89.17 | 76.91 | 88.88 |
| Zero-shot CoT | 55.10 | 73.14 | 91.89 | 80.68 | 91.60 | 91.84 | 74.52 | 58.59 | 65.03 | 85.35 | 89.12 | 75.91 | 89.22 |
| Zero-shot PoT | 18.46 | 6.35 | 48.09 | 71.66 | 82.20 | 83.73 | 54.20 | 32.40 | 48.71 | 78.86 | 86.74 | 67.18 | 59.49 |
| Zero-shot PS | 45.99 | 70.04 | 91.32 | 81.82 | 91.51 | 91.98 | 74.62 | 57.63 | 62.12 | 84.40 | 88.98 | 72.47 | 87.74 |
| Zero-shot CoT @maj16 | 69.18 | 82.68 | 93.56 | 85.35 | 92.56 | 94.51 | 83.92 | 66.46 | 74.86 | 89.46 | 92.27 | 84.45 | 92.70 |
| Zero-shot PoT @maj16 | 57.49 | 59.02 | 77.24 | 78.24 | 86.26 | 92.89 | 78.77 | 63.36 | 66.56 | 86.74 | 91.89 | 79.53 | 91.27 |
| Zero-shot PS @maj16 | 68.03 | 83.35 | 93.46 | 87.02 | 92.99 | 94.18 | 84.78 | 69.27 | 75.43 | 89.69 | 92.22 | 85.07 | 92.75 |
| PROVE | **74.90** | **84.83** | **93.76** | **87.88** | **93.35** | **94.75** | **87.02** | **76.67** | **80.49** | **89.98** | **92.83** | **86.35** | **93.03** |
| **MULTIARITH** | | | | | | | | | | | | | |
| Direct | 79.33 | 95.83 | 97.67 | 94.00 | 98.50 | 98.33 | 78.50 | 74.50 | 82.67 | 98.50 | 97.00 | 71.83 | 96.50 |
| Zero-shot CoT | 78.17 | 93.50 | 98.67 | 94.00 | 98.17 | 98.17 | 83.67 | 74.33 | 76.67 | 96.67 | 97.83 | 85.67 | 97.83 |
| Zero-shot PoT | 18.00 | 4.33 | 59.67 | 88.67 | 98.00 | 96.67 | 61.67 | 29.83 | 60.17 | 97.33 | 97.83 | 78.67 | 62.33 |
| Zero-shot PS | 62.33 | 87.50 | 97.83 | 94.33 | 98.67 | 98.50 | 83.50 | 70.67 | 70.00 | 97.83 | 97.33 | 83.67 | 97.00 |
| Zero-shot CoT @maj16 | 91.83 | 98.00 | **98.33** | 97.17 | 98.33 | **98.83** | 91.50 | 85.67 | 92.17 | 98.83 | 98.00 | 97.00 | **99.00** |
| Zero-shot PoT @maj16 | 78.33 | 49.67 | 90.33 | 95.00 | 98.67 | 98.67 | 92.83 | 77.83 | 79.00 | **99.00** | 99.00 | 95.83 | **99.00** |
| Zero-shot PS @maj16 | 91.33 | 98.33 | 98.17 | 97.00 | 98.67 | **98.83** | 94.33 | 90.33 | 92.17 | 98.17 | 98.50 | 97.83 | 98.50 |
| PROVE | **95.50** | **99.33** | **98.33** | **98.33** | **98.83** | **98.83** | **95.67** | **97.67** | **96.17** | **99.00** | **99.17** | **97.83** | **99.00** |
| **SINGLEEQ** | | | | | | | | | | | | | |
| Direct | 72.44 | 88.19 | 95.47 | 91.14 | 94.88 | 96.06 | 85.24 | 79.92 | 82.68 | 93.50 | 93.31 | 89.96 | 94.69 |
| Zero-shot CoT | 67.72 | 85.24 | 97.44 | 92.32 | 96.65 | 98.23 | 83.86 | 78.54 | 83.07 | 93.70 | 94.49 | 87.40 | 95.47 |
| Zero-shot PoT | 25.00 | 5.91 | 49.61 | 94.49 | 98.23 | 98.03 | 69.49 | 41.34 | 69.69 | 94.49 | 96.26 | 84.06 | 62.40 |
| Zero-shot PS | 62.80 | 81.69 | 96.65 | 92.13 | 97.24 | 98.03 | 85.83 | 75.00 | 76.57 | 91.14 | 95.28 | 84.45 | 93.90 |
| Zero-shot CoT @maj16 | 81.50 | 90.75 | 96.85 | 94.49 | 97.24 | 98.82 | 89.37 | 83.27 | 92.13 | 96.85 | 97.24 | 94.09 | 97.83 |
| Zero-shot PoT @maj16 | 80.31 | 71.06 | 79.33 | 96.06 | 98.82 | 99.02 | **94.69** | 86.22 | 88.78 | 98.03 | 98.43 | 94.49 | 97.64 |
| Zero-shot PS @maj16 | 80.71 | 93.11 | 97.44 | 95.08 | 98.03 | 98.23 | 91.14 | 85.24 | 91.73 | 97.24 | 97.83 | 94.88 | 97.24 |
| PROVE | **92.13** | **96.80** | **98.82** | **97.64** | **99.21** | **99.21** | **94.69** | **92.91** | **97.24** | **98.23** | **98.62** | **96.06** | **98.82** |
| **SINGLEOP** | | | | | | | | | | | | | |
| Direct | 77.76 | 90.75 | 96.09 | 91.28 | 95.37 | 96.8 | 90.04 | 83.45 | 88.79 | 93.95 | 95.73 | 92.35 | 96.26 |
| Zero-shot CoT | 6.33 | 89.32 | 96.26 | 91.64 | 95.91 | 96.26 | 90.93 | 85.41 | 85.94 | 94.13 | 96.09 | 88.43 | 95.91 |
| Zero-shot PoT | 29.54 | 10.5 | 44.84 | 90.57 | 95.73 | 95.55 | 72.78 | 51.25 | 68.51 | 93.24 | 95.91 | 84.70 | 74.02 |
| Zero-shot PS | 66.19 | 88.26 | 96.26 | 93.42 | 96.26 | 96.80 | 90.75 | 81.67 | 81.32 | 91.10 | 95.55 | 86.83 | 94.13 |
| Zero-shot CoT @maj16 | 87.37 | 92.88 | 96.26 | 94.13 | 96.09 | 96.80 | 94.66 | 88.79 | 92.70 | 96.62 | 96.80 | 93.77 | **96.80** |
| Zero-shot PoT @maj16 | 84.52 | 77.22 | 80.60 | 92.88 | 95.73 | 96.80 | 93.59 | 89.86 | 85.77 | 95.37 | 96.44 | 91.99 | 96.44 |
| Zero-shot PS @maj16 | 88.08 | 93.59 | 96.26 | 94.66 | **96.62** | 97.15 | 94.84 | 91.10 | 91.28 | 96.44 | 96.62 | 93.06 | 96.62 |
| PROVE | **92.53** | **94.66** | **96.98** | **94.84** | **96.62** | **97.51** | **95.02** | **94.66** | **93.24** | **97.16** | **96.98** | **94.66** | **96.80** |
| **ADDSUB** | | | | | | | | | | | | | |
| Direct | 61.27 | 80.51 | 92.91 | 84.05 | 90.38 | 93.16 | 81.01 | 73.16 | 82.53 | 88.10 | 89.62 | 84.56 | 90.38 |
| Zero-shot CoT | 58.73 | 79.75 | 94.94 | 86.33 | 94.94 | 95.95 | 82.28 | 73.67 | 81.01 | 89.37 | 92.66 | 84.56 | 93.42 |
| Zero-shot PoT | 49.62 | 7.34 | 28.61 | 53.92 | 95.95 | 95.7 | 70.38 | 42.53 | 65.57 | 93.67 | 95.44 | 78.73 | 88.86 |
| Zero-shot PS | 51.39 | 77.47 | 92.91 | 89.11 | 95.19 | 94.43 | 83.29 | 73.42 | 75.19 | 85.32 | 92.66 | 77.97 | 90.89 |
| Zero-shot CoT @maj16 | 70.89 | 86.58 | 95.44 | 90.89 | 94.43 | 96.2 | 85.82 | 80.51 | 86.08 | 93.42 | 95.19 | 89.62 | **95.95** |
| Zero-shot PoT @maj16 | 71.39 | 64.81 | 63.54 | 90.89 | 95.95 | 96.46 | 94.11 | 81.52 | 85.06 | **95.19** | 95.95 | 82.53 | **95.95** |
| Zero-shot PS @maj16 | 71.14 | 85.82 | 94.94 | 91.14 | 95.7 | 96.46 | 88.35 | 83.29 | 87.34 | 93.16 | 94.43 | 90.13 | 94.68 |
| PROVE | **85.06** | **92.15** | **96.46** | **92.41** | **96.20** | **96.96** | **91.65** | **89.37** | **93.16** | **95.19** | **96.20** | **92.41** | **95.95** |

Table 1: Main results showing the comparison of PROVE and baseline methods across 13 LLMs, ranging in size from 0.5B to 13B, on seven mathematical reasoning datasets.
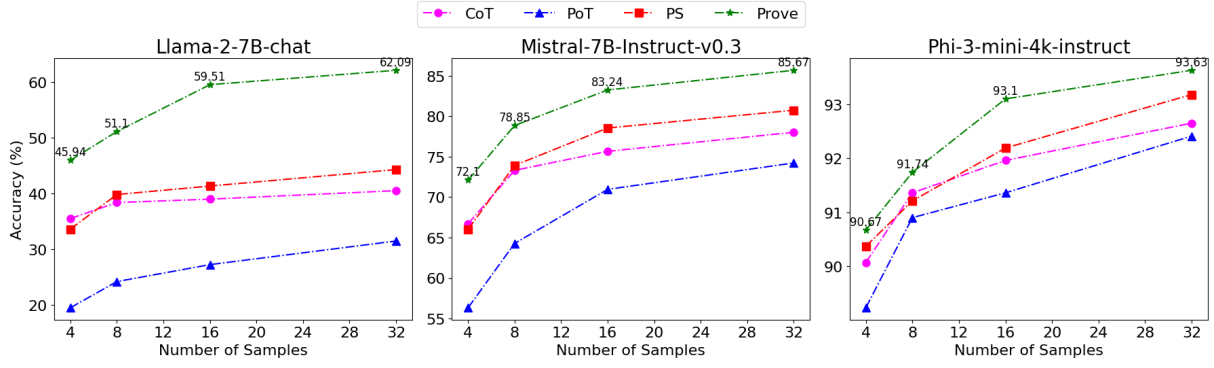
Figure 5: Performance comparison of PROVE and baseline methods across various LLMs evaluated on GSM8K with different number of samples.

## 3.2 Main Results

We report the main results of our experiments in Table 1. We highlight several key observations:

**PROVE improves over baselines consistently on all evaluated LLMs across seven datasets.** Our evaluations of PROVE on various model families and sizes reveal consistent improvements over strong baseline models. For instance, on the GSM8K dataset, Mistral 7B achieves an accuracy of 83.24% with PROVE, compared to 78.54% with the best baseline (Zero-shot PS @maj16), reflecting an improvement of 4.7%. In the SVAMP dataset, Mistral 7B reaches 88.90%, surpassing the best baseline by 4.2%. Similarly, in the ASDIV dataset, Mistral 7B achieves an accuracy of 87.02%, exceeding the baseline's 84.78% by 2.24%. Furthermore, Mistral 7B shows consistent performance gains of up to 3.39% in other datasets, including MultiArith, SingleEq, SingleOp, and AddSub. These improvements are consistent across all twelve LLMs we evaluated.

**PROVE provides larger performance gains for smaller models.** In our comparison of various model families, including Qwen 2, Gemma 2, Llama 2, and Llama 3.2, we found that smaller models tend to exhibit greater performance gains than their larger counterparts. For example, in the GSM8K dataset, the smaller Llama 2 7B model shows an improvement of 18.19% over the baseline, while the larger Llama 2 13B model only achieves a gain of 13.87%. Similarly, Gemma 2 7B demonstrates a 6.22% increase in accuracy compared to the baseline, whereas the larger Gemma 2 9B model achieves only a 0.76% improvement. The Qwen 2 models also reflect a similar trend, the 0.5B variant outperforms its baseline by 4.98%, while

the 1.5B model improves by 2.65%. In contrast, the larger Qwen 2 7B model shows only a 0.76% improvement. Likewise, Llama 3.2 1B demonstrates a 7.35% improvement compared to its larger variant, the 3B model, which only achieves a 2.5% increase. These results indicate that smaller models benefit more significantly from PROVE, likely because they are more prone to errors, such as miscalculations, making the filtering more beneficial for them.

**Newer models closing the gap with larger ones.** More recent Llama models, such as Llama 3.2, demonstrate significant progress in closing the performance gap with older, larger models like Llama 2 and Llama 3, despite their smaller size. For example with PROVE on datasets like GSM8K, Llama 3.2 1B outperforms Llama 2 13B achieving 73.01% accuracy versus 68.08%, and Llama 3.2 3B outperforms Llama 3 8B achieving 91.36% over 90.14%. On datasets like ASDIV, the performance of Llama 3.2 1B and Llama 3.2 3B increases the separation even further, reflecting improvements in architecture and training techniques to increase performance of smaller models (Meta, 2024).

## 3.3 Ablation Study

**Impact of the Number of Plans.** In Figure 5, we compare the accuracy of CoT, PoT, and PS prompting methods with SC and PROVE across different numbers of sampled solutions on GSM8K. As the number of sampled solutions increases from 4 to 32, PROVE consistently outperforms the baseline methods, demonstrating its effectiveness. Notably, with just 4 samples, the Llama-2-7b-chat and Mistral-7B-Instruct-v0.3 models show significant accuracy gains of 10.46% and 5.39%, respectively. With 32 samples, Llama-2-7b-chat and Mistral-7B-Instruct-v0.3 models have a gain of 17.81% and 4.93%.

**Question:**
A mother goes shopping. She buys cocoa at $4.20, laundry at $9.45 and a package of pasta at $1.35. She pays $20. How much change does the cashier give back?

**Generated Solution:**
To find the change the cashier gives back, we need to subtract the total cost of the items purchased from the amount the mother paid.

Step 1: Add the costs of all items purchased:
- Cocoa: $4.20
- Laundry: $9.45
- Pasta: $1.35
Total cost = $4.20 + $9.45 + $1.35 = $15.10

Step 2: Subtract the total cost from the amount paid:
Amount paid = $20
Change = Amount paid - Total cost
Change = $20 - $15.10
Change = $4.90

Therefore, the cashier gives back $4.90 as change.

**Translated Python Program:**
```python
def solution():
    # Step 1: Add the costs of all the items purchased
    cocoa_cost = 4.20
    laundry_cost = 9.45
    pasta_cost = 1.35
    total_cost = cocoa_cost + laundry_cost + pasta_cost

    # Step 2: Subtract the total cost from the amount paid
    amount_paid = 20
    change = amount_paid - total_cost

    # Return the change
    return change
```

**Question:**
Jane counts two zebras with 17 stripes each, a zebra with 36 stripes, and another zebra with half that many stripes. How many stripes do the zebras have on average?

**Generated Solution:**
To find the average number of stripes on the zebras, we need to calculate the total number of stripes and then divide it by the number of zebras.

Step 1: Calculate the total number of stripes.
- Jane's zebras have 17 stripes each, so she has 2 * 17 = 34 stripes.
- The zebra with 36 stripes contributes 36 stripes.
- The zebra with half the number of stripes (36 / 2 = 18) contributes 18 stripes.
Total number of stripes = 34 + 36 + 18 = 88

Step 2: Calculate the average number of stripes.
There are 3 zebras in total. To find the average, we divide the total number of stripes by the number of zebras:
Average number of stripes = Total number of stripes / Number of zebras
= 88 / 3 = 29.33

So, the average number of stripes on the zebras is approximately 29.33.

**Translated Python Program:**
```python
def solution():
    # Step 1: Calculate the total number of stripes
    zebra_1_stripes = 17
    zebra_2_stripes = 17
    zebra_3_stripes = 36
    zebra_4_stripes = zebra_3_stripes / 2
    total_stripes = (zebra_1_stripes * 2) + zebra_3_stripes + zebra_4_stripes

    # Step 2: Calculate the average number of stripes per zebra
    num_zebras = 4
    average_stripes = total_stripes / num_zebras

    return average_stripes
```

Figure 6: Case study examples where PROVE gets correct when vanilla majority voting fails. The generated solution is by Llama-3.2-1B-Instruct on a GSM8K type question. In both examples, we show how the translated program can correctly filter out generated solution that have errors in them. Left: highlighted in red shows a calculation error. Right: highlighted in red shows semantic understanding error.

| Model | Translation Model | Accuracy (%) |
|---|---|---|
| Llama-3.2-1B-Instruct | Phi-3-mini-4k-instruct | 73.01 |
| | Llama-3-8B-Instruct | 73.54 |
| | Llama-3.1-8B-Instruct | 74.53 |
| Llama-2-7B-chat | Phi-3-mini-4k-instruct | 59.51 |
| | Llama-3-8B-Instruct | 58.15 |
| | Llama-3.1-8B-Instruct | 61.94 |
| Mistral-7B-Instruct-v0.3 | Phi-3-mini-4k-instruct | 83.24 |
| | Llama-3-8B-Instruct | 80.29 |
| | Llama-3.1-8B-Instruct | 83.32 |

Table 2: Performance of PROVE on GSM8K, evaluated across various translation models and LLMs.

| Model | Prompting Method | Accuracy (%) |
|---|---|---|
| Llama-3.2-1B-Instruct | Direct | 75.59 |
| | CoT | 75.89 |
| | PS | 73.01 |
| Llama-2-7B-chat | Direct | 59.06 |
| | CoT | 59.29 |
| | PS | 59.51 |
| Mistral-7B-Instruct-v0.3 | Direct | 82.34 |
| | CoT | 82.79 |
| | PS | 83.24 |

Table 3: Performance of PROVE on GSM8K, with different prompting strategies across various LLMs.

**Impact of Translation Model.** We study the impact of model selection for translating plans and solutions into Python programs on GSM8K. In our experiments, the default translation model used is Phi-3-mini-4k-instruct. We tested Llama-3-8B-Instruct and Llama-3.1-8B-Instruct for translating the outputs of LLMs, which include Llama-3.2-1B-Instruct, Llama-2-7B-chat, and Mistral-7B-Instruct-v0.3. Table 2 shows the performance variations across different translation models for the different LLMs. We observe that Llama-3.1-8B-Instruct as the translation model, with more parameters than Phi-3-mini-4k-instruct, leads to slight performance improvements for all tested LLMs, having performance gains of 1.52%, 2.43%, and 0.08% for Llama-3.2-1B-Instruct, Llama-2-7B-chat, and Mistral-7B-Instruct-v0.3, respectively.

**Impact of Plan and Solve (PS) Prompting.** We investigated the impact of PS prompting on our framework by comparing its effectiveness against both direct prompting and CoT prompting in the first stage of PROVE. We evaluated on three different models, namely Llama-3.2-1B-Instruct, Llama-2-7B-Chat, and Mistral-7B-Instruct-v0.3, on GSM8K. As shown in Table 3, the overall performance differences between the prompting methods are relatively small. The largest difference observed is 2.88%, with CoT outperforming PS when evaluating on Llama-3.2-1B-Instruct. Our findings indicate that even with variations in the prompting method during the first stage, PROVE consistently maintains performance, demonstrating that it is prompt-agnostic in this stage of the framework.

## 3.4 Qualitative Analysis

**PROVE gets correct when vanilla majority vote fail.** In Figure 6, we show two case study examples of PROVE that correctly filter generated solutions that have errors in them. These two examples are examples that majority voting fails to solve due to the correct answer not being the most common answer, whereas, PROVE is able to successfully filter out incorrect reasoning paths, making the correct answer as the majority answer. In Figure 6 (left), the model did a wrong calculation "$\$4.20 + \$9.45 + \$1.35 = \$15.10$", whereas the correct calculation should be "$\$15.00$". However in our PROVE framework, the translated program is able to detect the miscalculation from the generated solution and filter it out. Similarly in Figure 6 (right), the model had a semantic understanding error where it fails to understand that the total number of zebra should be 4 instead of 3. Again, our translated program is able to point out the correct value and filter it out this generated solution. These examples demonstrate that by leveraging programs as verifiers, we can effectively filter out reasoning paths that contain errors, such as miscalculations and semantic misunderstandings.

| Method | Calculation | Missing | Semantic |
|---|---|---|---|
| Zero-shot PS @maj16 | 20% | 12% | 68% |
| PROVE | 0% | 18% | 82% |

Table 4: Distribution of error types (calculation error, missing step error, semantic understanding error) for 50 randomly sampled examples from GSM8K, where both PS and PROVE got incorrect, using Llama-3.2-1B-Instruct.

**Error Analysis.** To better understand the effectiveness of PROVE, we conducted a manual analysis of 50 randomly sampled examples from GSM8K, where both PS and PROVE got incorrect, using Llama-3.2-1B-Instruct. Following (Wang et al., 2023a), we categorized the errors into three types: calculation errors, missing step errors, and semantic understanding errors. The results of this analysis (Table 4) show that PROVE achieves a 0% calculation error rate, significantly lower than PS, which had a 20% calculation error rate. This demonstrates that PROVE effectively minimizes calculation errors, as illustrated in the case study example shown in Figure 6 (left). Additionally, in some instances, it also reduces semantic understanding errors, as depicted in Figure 6 (right).

## 4 Related Works

**Reasoning with Large Language Models.** Recent advances in LLMs have demonstrated remarkable capabilities in tackling complex reasoning tasks. Prior research highlights that step-by-step reasoning prompts, such as CoT (Kojima et al., 2022) and PS (Wang et al., 2023a), improve performance compared to directly generating answers. Other techniques, such as multi-step decoding (Yao et al., 2023), explore diverse reasoning paths, while methods like multi-sample reasoning (e.g., majority voting (Wang et al., 2023b)) aggregate solutions to enhance robustness. However, these approaches still heavily depend on the LLM's generative outputs, which are prone to hallucination, particularly in smaller models. Our method, PROVE, complements these by introducing a verification step that uses Python programs to filter out potentially flawed reasoning paths, resulting in more accurate solution aggregation.

**Large Language Models as Verifiers.** Using language models to evaluate model generations has been a long standing idea(Kushman et al., 2014; Roy and Roth, 2015; Shen et al., 2021). A common approach involves training a separate verification model to assess the correctness of generated outputs (Cobbe et al., 2021). Other techniques, such as step-by-step verification (Lightman et al., 2024) and ranking multiple reasoning paths to choose the most accurate one (Weng et al., 2023), show promise in reducing errors but often rely heavily on the inherent reasoning abilities of the model, limiting their effectiveness, especially in smaller models. While program-based verification approaches have been proposed (Zhou et al., 2024a; Han et al., 2024; Zhou et al., 2024b), they tend to be more complex, often requiring few-shot prompting and more capable models. In contrast, PROVE is simpler, does not require few-shot examplars, and is easily adaptable to smaller, open-source LLMs.

## 5 Conclusion

In this paper, we demonstrate how using programs as verifiers can effectively filter out reasoning paths that contain calculation or semantic understanding errors. Our approach, PROVE, is model-agnostic and does not require fine-tuning or few-shot examplars for prompting. PROVE consistently outperform baseline methods across 13 LLMs and seven mathematical reasoning datasets.

# 6 Limitations

One limitation of this approach is that Python programs cannot directly parse LaTeX syntax, making it challenging to evaluate datasets that heavily rely on LaTeX, such as MATH (Hendrycks et al., 2021). As a result, due to the simplicity of our method, PROVE is currently limited to working only on math datasets where the final answer is numeric. Future work will focus on adapting our framework to accommodate LaTeX-intensive math datasets.

# References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Qin Cai, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Yen-Chun Chen, Yi-Ling Chen, Parul Chopra, Xiyang Dai, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Victor Fragoso, Dan Iter, Mei Gao, Min Gao, Jianfeng Gao, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Ce Liu, Mengchen Liu, Weishung Liu, Eric Lin, Zeqi Lin, Chong Luo, Piyush Madan, Matt Mazzola, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Xin Wang, Lijuan Wang, Chunyu Wang, Yu Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Haiping Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Sonali Yadav, Fan Yang, Jianwei Yang, Ziyi Yang, Yifan Yang, Donghan Yu, Lu Yuan, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *Preprint*, arXiv:2404.14219.

Sourav Banerjee, Ayushi Agarwal, and Saloni Singla. 2024. Llms will always hallucinate, and we need to live with this. *Preprint*, arXiv:2409.05746.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gu-

rurangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2024. ToRA: A tool-integrated reasoning agent for mathematical problem solving. In *The Twelfth International Conference on Learning Representations*.

Vernon Toh Yan Han, Ratish Puduppully, and Nancy F. Chen. 2024. Veritymath: Advancing mathematical reasoning by self-verification through unit consistency. *Preprint*, arXiv:2311.07172.

Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8154–8173, Singapore. Association for Computational Linguistics.

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. *NeurIPS*.

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.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*.

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.

Nate Kushman, Yoav Artzi, Luke Zettlemoyer, and Regina Barzilay. 2014. Learning to automatically solve algebra word problems. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 271–281, Baltimore, Maryland. Association for Computational Linguistics.

Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*.

Lei Lin, Jiayi Fu, Pengli Liu, Qingyang Li, Yan Gong, Junchen Wan, Fuzheng Zhang, Zhongyuan Wang, Di Zhang, and Kun Gai. 2024. Just ask one more time! self-agreement improves reasoning of language models in (almost) all scenarios. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 3829–3852, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. In *Advances in Neural Information Processing Systems*, volume 36, pages 46534–46594. Curran Associates, Inc.

Meta. 2024. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models.

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.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex

Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

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.

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.

Jianhao Shen, Yichun Yin, Lin Li, Lifeng Shang, Xin Jiang, Ming Zhang, and Qun Liu. 2021. Generate & rank: A multi-task framework for math word problems. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2269–2279, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. 2024. Gemma 2: Improving open language models at a practical size. *Preprint*, arXiv:2408.00118.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2024. Mathcoder: Seamless code integration in LLMs for enhanced mathematical reasoning. In *The Twelfth International Conference on Learning Representations*.

Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023a. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2609–2634, Toronto, Canada. Association for Computational Linguistics.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023b. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.

Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. 2023. Large language models are better reasoners with self-verification. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2550–2575, Singapore. Association for Computational Linguistics.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan,

Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. Qwen2 technical report. *Preprint*, arXiv:2407.10671.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. In *Thirty-seventh Conference on Neural Information Processing Systems*.

Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2024a. Solving challenging math word problems using GPT-4 code interpreter with code-based self-verification. In *The Twelfth International Conference on Learning Representations*.

Jin Peng Zhou, Charles E Staats, Wenda Li, Christian Szegedy, Kilian Q Weinberger, and Yuhuai Wu. 2024b. Don't trust: Verify – grounding LLM quantitative reasoning with autoformalization. In *The Twelfth International Conference on Learning Representations*.