

# Interpretable Math Word Problem Solution Generation Via Step-by-step Planning

Mengxue Zhang<sup>1</sup>, Zichao Wang<sup>2</sup>, Zhichao Yang<sup>1</sup>, Weiqi Feng<sup>1</sup>, Andrew Lan<sup>1</sup>

<sup>1</sup>University of Massachusetts Amherst, <sup>2</sup>Rice University

mengxuezhang@umass.edu

## Abstract

Solutions to math word problems (MWP) with step-by-step explanations are valuable, especially in education, to help students better comprehend problem-solving strategies. Most existing approaches only focus on obtaining the final correct answer. A few recent approaches leverage intermediate solution steps to improve final answer correctness but often cannot generate coherent steps with a clear solution strategy. Contrary to existing work, we focus on improving the correctness and coherence of the intermediate solutions steps. We propose a step-by-step planning approach for intermediate solution generation, which strategically plans the generation of the next solution step based on the MWP and the previous solution steps. Our approach first *plans* the next step by predicting the necessary math operation needed to proceed, given history steps, then *generates* the next step, token-by-token, by prompting a language model with the predicted math operation. Experiments on the GSM8K dataset demonstrate that our approach improves the accuracy and interpretability of the solution on both automatic metrics and human evaluation.

## 1 Introduction

Arithmetic math word problems (MWPs) consist of natural language statements describing real-world scenarios that involve numerical quantities, followed by a question asking for an unknown value. Solving MWPs require parsing the textual statements and carrying out the corresponding calculations (Kumar et al., 2022). MWPs are an important educational tool that helps assess and improve student knowledge in basic mathematical concepts and skills (Walkington, 2013; Verschaffel et al., 2020). They also represent a long-standing interest in artificial intelligence (AI) research since correctly solving them serves as a key benchmark task for testing and improving the mathematical reasoning skills of AI models (Feigenbaum and Feldman,

1995; Bommasani et al., 2021; Cobbe et al., 2021; Lewkowycz et al., 2022).

There is a large body of literature that focuses on automatically solving MWP. Earlier works took a modular approach that first analyzes unconstrained natural language and then maps intricate text patterns onto mathematical vocabulary (Sundaram et al., 2022). As a result, this approach relies heavily on hand-crafted rules to fill the gap between natural language and symbolic mathematical vocabulary (Sundaram et al., 2022). Recent works leverage advances in natural language processing and take a neural network-based, end-to-end approach, where a neural network encodes a numerical representation of the MWP (and the underlying equation), from which a decoder generates the final answer (Zou and Lu, 2019; Wang et al., 2017; Wu et al., 2020; Chen et al., 2020; Cao et al., 2021; Shen et al., 2021; Shao et al., 2022; Jie et al., 2022). Unfortunately, the vast majority of these works focus on generating and predicting *a single final answer*, since answer correctness is often the only evaluation metric. Therefore, these works do not provide any insights or explanations into how the models arrive at the answer. As a result, it is often difficult, if not entirely impossible, to explain the model’s behavior, especially when it produces a wrong answer. The lack of interpretability of these methods makes it challenging to analyze them and unsafe to use them in real-world applications.

This interpretability issue has attracted increasing interest in MWP solving research. Recent works have shifted to designing models that not only generate the final answer for an MWP, *but also the intermediate steps*. The ability to generate intermediate steps not only enables researchers to investigate model behavior but also new applications. For example, in personalized education and intelligent tutoring systems, these models have the potential to generate detailed, personalized solution steps as feedback to improve stu-

dent understanding of the mathematical concepts and resolve misconceptions (Walkington, 2013; Karpicke, 2012; Koedinger et al., 2015). The recent GSM8K (Cobbe et al., 2021) dataset contains MWP that come with 2 to 8 intermediate steps described in natural language, which provides us a good resource to study step-by-step solution generation. Many works apply (large) language models (LMs) on this dataset and achieve high accuracy in final answer generation, without studying the quality of intermediate steps (Wei et al., 2022; Wang et al., 2022; Chowdhery, Aakanksha and others, 2022; Lewkowycz et al., 2022; Uesato et al., 2022; Kojima et al., 2022; Li et al., 2022). These works use verifiers, self-consistency decoding strategy (majority votes), chain-of-thought prompting, or calculators; see Section 4 for a detailed discussion.

However, existing LMs are still prone to generating incorrect intermediate steps despite yielding the correct final answer. The models are not competent at numerical reasoning, possibly because they generate intermediate steps word by word (or token by token) and cannot look far ahead. As a result, they only use shallow heuristics (Li et al., 2021) in word occurrence and lack multi-step mathematical reasoning capabilities, which solving an MWP requires. A recent study that experiments on GPT-4 also points out that the architecture of next-word prediction precludes any “inner dialog” and cannot really plan ahead (Bubeck et al., 2023).

## 1.1 Contributions

In this paper, we study the problem of generating accurate and high-quality intermediate solution steps with natural language explanation via step-by-step planning using LMs. We formulate this problem as a controllable generation problem where the LM aims to generate the correct intermediate solution at each solution step, given the MWP and previous solution steps. This problem is particularly challenging since *the generated solution steps need to be accurate*, i.e., each intermediate step must be mathematically valid and on the path to the correct answer. We need an approach different from widely-adopted, attribute-controlled generation approaches for topic or sentiment, where the attribute is nuanced and cannot be matched exactly (Dathathri et al., 2020; Krause et al., 2020; Shirish Keskar et al., 2019).

To overcome these challenges, we introduce a *planning-LM* approach, where we plan the strategy

for the next solution step and then use the plan to guide LMs to generate the step. Since symbols and patterns are crucial to the effectiveness of chain-of-thought prompting (Madaan and Yazdanbakhsh, 2022), we design plans in the form of *mathematical operations* to prompt the model to generate the next intermediate step. We summarize our contributions as follows.

[C1] We explore the use of a planning approach for step-by-step solution generation for MWPs. To the best of our knowledge, our work is the first to focus on generating high-quality intermediate solution steps via LMs.

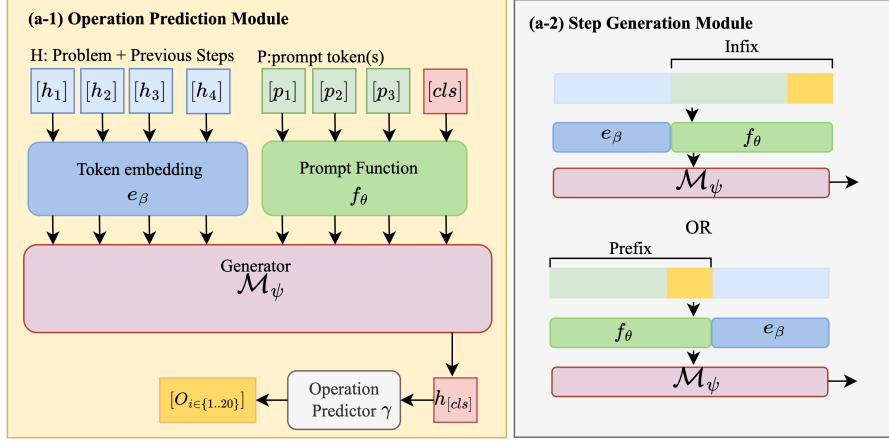
[C2] We first predict the mathematical operation applied in the next solution step using a small model and then apply a carefully-constructed prompt to control an LM to generate the next solution step. Our approach can be extended to many downstream applications due to its interpretability and high controllability.

[C3] We evaluate our planning-LM approach on the GSM8K dataset to demonstrate its effectiveness, both quantitatively and qualitatively. With minimal additional parameters (0.02%), it outperforms existing approaches on both final answer accuracy and intermediate step quality. Moreover, by manually changing the math operation prompt, we can control our approach to generate *different correct solution paths* for the same MWP.

## 1.2 Notation

We first define all of the terms and components in our approach. We define an MWP as  $Q = \{q_1, q_2, \dots, q_n\}$  where  $q_i$  represents a token, which is either a numerical value, a mathematical operator, or a word/sub-word. The corresponding step-by-step solution is  $S = \{S^1, S^2, \dots\}$ , where  $S^i$  denotes  $i^{th}$  step of the solution. For any step  $S^i$ , we denote it as  $S^i = \{s_1^i, s_2^i, \dots\}$ , consisting of a sequence of tokens. Next, we define our prompt in two parts. The first part is the textual instruction prompt, which contains words that LMs can understand, and the second part is the mathematical operation prompt, which is a special token that instructs the LM to perform which mathematical operation in the next solution step. We denote the instruction prompt as  $P = \{p_1, p_2, \dots\}$ , where  $p_i$  represents a word/sub-word token, and the operation prompt as  $O = \{o\}$ , where  $o$  is a categorical variable indicating the math operation token. We define  $H_i$  as the solution context, i.e., the history

(a) Basic Components



(b) Generation via step-by-step planning

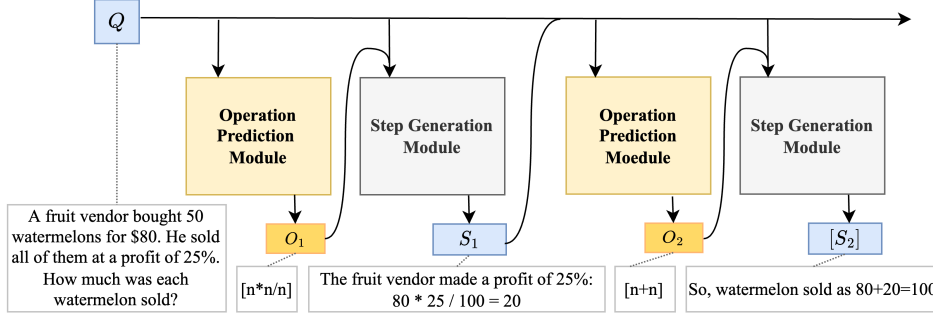


Figure 1: An overview of our step-by-step MWP solution generation approach. Planning-LM first predicts the next step operation hint (a-1) and controls the next step generate via the predicted operation hint (a-2). **Figure (b)** shows the overview generation process by given the question  $Q$ .

at step  $S^i$ , which consists of the problem  $Q$  and all previous steps,  $\{S^1, \dots, S^{i-1}\}$ .  $\mathcal{M}$  denotes the base LM and  $e$  is its corresponding token embedding function. Finally, we define  $f$  as the prompt embedding function. Both  $e$  and  $f$  can map tokens into  $\mathcal{R}^K$  where  $K$  is the hidden state dimension of the LM.

## 2 Methodology

We now define our MWP solution generation task and detail the specifics of our approach. Our task is that given a question  $Q$ , we need to generate a step-by-step solution  $S = S^1, S^2, \dots$ , with each step consisting of a combination of textual and mathematical tokens, to reach the final answer. We formulate the problem as a step-wise controllable generation task using prompts-based LM fine-tuning. Figure 1 shows an overview of our approach<sup>1</sup> in-

<sup>1</sup>For clarity, we discuss our methodology based on decoder-only Transformer-based LMs. However, our methodology also generalizes to encoder-decoder-type LMs, such as T5, which we experimentally verify (see Table 1). More details can be found in Appendix E.

cluding its two main components: First, we utilize the MWP and the solution history to plan and predict next mathematical operation to apply in the next step. Second, we use the predicted operation prompt with instruction prompt to guide the next step generation process. Our key technical challenges are (i) how to learn a solution planning strategy to transition from step to step and (ii) once we have the next operation, how to apply and design prompts to guide the generative LM to generate the next step to follow the plan.

### 2.1 Operation Prediction

Our first step is to predict the mathematical operation to be applied in the next step. To achieve this, we concatenate the solution history  $H$  and a crafted instruction prompt  $P$  (e.g., "What is the next operation?") followed by the special token  $[cls]$  as input to an (not necessarily large) LM. We encode solution history tokens with a vocabulary embedding function  $e_\beta$  and instruction prompt tokens with a separate prompt embedding function

$f_\theta$ ;  $\beta$  and  $\theta$  are the parameters of these parts, i.e., the embedding layer in an LM. Then, we obtain the representation of the solution history as the final layer hidden state of the LM, i.e.,  $\mathcal{M}$ . To predict the operation action of the next step, we use a one-layer, fully-connected network as the classifier, with weight  $w_\gamma$ , to obtain an operation score vector for each valid math operation  $s \in [0, 1]^{|O|}$ , where  $|O|$  is the number of operation classes, as

$$s = w_\gamma \dot{h}_{[cls]},$$

where  $\gamma$  is the set of parameters for the classifier. Since we need to use an LM for step generation, introducing a separate LM for operation prediction leads to a large number of parameters. Therefore, we use the same LM for both operation planning and solution step generation. The objective function for operation planning is the cross-entropy loss on operators, i.e.,

$$\mathcal{L}_{CE} = - \sum_i^{|O|} t_i \log \left( \frac{\exp s_i}{\sum_j^{|O|} \exp s_j} \right),$$

where  $s_i$  is the score of operation class  $i$ .  $t_i$  is an indicator such that  $t_i = 1$  when  $i$  is the true label and  $t_i = 0$  otherwise. We obtain true labels by extracting mathematical operations from each step of the solution in the training data, which we detail below in Section 2.3.

## 2.2 Controllable Step Generation

Once we have the predicted operation  $O$ , we append the corresponding prompt to the instruction prompt  $P$  to form our final prompt for step generation. Our task becomes a controllable generation task: given history  $H$  and the prompt  $[P; O]$  that plans the next step, our goal is to generate the next step  $S$  token-by-token. We generate a step  $S_i = \{s^i_1, \dots, s^i_T\} = \{s^i_j\}_{j=1}^T$  according to

$$p(S_i | [P_i; O_i], H_i) = \prod_{j=1}^T p(s^i_j | [P_i; O_i], H_i, \{s^i_j\}_{j=1}^{j-1}).$$

Then, the overall step-by-step solution  $S$  with  $N$  steps is generated according to

$$p(S) = \prod_{i=1}^N p(S_i | [P_i; O_i], H_i) p(O_i | H_i).$$

The step generation objective is given by the negative log-likelihood objective function

$$\mathcal{L}_{LM} = - \sum_{i=1}^N \log p_{\beta, \theta, \gamma, \psi}(S_i | [P_i; O_i], H_i),$$

where the set of parameters include previously defined  $\beta, \theta, \gamma$  and the LM parameters  $\psi$ .  $\beta$  and  $\psi$  are fine-tuned while  $\theta$  and  $\gamma$  are learned from scratch. We also investigate two ways to position the prompt in LM input: as prefix, where we place them at the beginning, i.e., the input is given by  $[P; O; H]$  and as infix, where we append the prompt after the history, i.e., the input is given by  $[H; P; O]$ .

## 2.3 Prompt Design

Our prompt consists of two parts: the instruction prompt gives the LM general instructions on what to generate, while the operation prompt provides specific guidelines for the mathematical calculation involved in the next step. For the instruction prompt, we apply prompt mining (Yuan et al., 2021) to find good instructions, i.e., word tokens that are the most informative for the LM to accomplish the desired task. See Section D.2 for details. For the operation prompt, we extract 20 common operations from the training data, such as one step addition  $[n + n]$ , subtraction  $[n - n]$ , multiplication  $[n * n]$ , etc and use them as prompts. We note that these operators are easy to find and can be automatically extracted, which means that there is no need to manually create labels to train the operation prediction LM. The instruction tokens and operation action tokens form the entire vocabulary of the prompt function  $f_\theta$ . The prompt function is a two-layer perceptron with a ReLU activation function.

## 2.4 Optimization

Although our entire approach can be trained together in an end-to-end way, we found that optimizing the operation prediction model and fine-tuning the LM/prompts for step generation asynchronously leads to better performance. Our intuition is that the operation predictor is a high-level decision-making policy for the entire solution while the LM generation process is a low-level (token-by-token) decision-making process for the current step. Optimizing these two modules simultaneously may cause inconsistency since the operation predictor may make a decision based on LM parameters that also need to be updated. Therefore, we first optimize the parameters of the generation LM and prompts with the step generation task loss, using ground truth operation labels, which we extract from the mathematical part of each step in the training data. Then, we iterate between freezing both the LM  $M$  and the prompt function  $f$  while



Table 1: Planning-LM outperforms fine-tuning LMs for both small and medium-sized GPT-2. Moreover, Planning-LM with a small GPT-2 achieves performance comparable to fine-tuning medium GPT-2, implying that our method can make a smaller model rival larger ones fine-tuned in the traditional way.

Model	BLEU	ACC-eq	ACC-op	Solve Rate
Chain-of-thought-tuning GPT-2 (117M)	34.3	49.4	55.1	8.1
Planning-GPT-2 with operation classifier (117M)	35.4	56.7	61.6	14.1
Chain-of-thought-tuning GPT-2-medium (345M)	38.1	58.1	61.1	16.1
Planning-GPT-2-medium with operation classifier (345M)	39.5	61.8	65.2	20.1
Chain-of-thought-tuning T5 (220M)	30.3	45.4	52.1	3.1
Planning-T5 with operation classifier (220M)	34.4	55.7	60.6	13.9
Chain-of-thought-tuning T5-large (770M)	35.3	58.9	63.1	17.0
Planning-T5-large with operation classifier (770M)	40.5	62.3	66.3	21.2

Table 2: Ablation results for different components of our approach. Most components contribute significantly.

Method Component					Metric			
Infix	Prefix	Prompt function	Prompt mining	Operation Predictor	BLEU	ACC-eq	ACC-op	Solve Rate
✓		✓	✓	✓	<b>35.4</b>	<b>56.7</b>	61.6	<b>14.1</b>
	✓	✓	✓	✓	33.7	52.1	<b>63.2</b>	10.4
✓			✓	✓	33.1	51.9	58.4	10.2
✓		✓		✓	33.9	55.1	59.9	13.2
✓		✓	✓		34.1	54.2	60.1	13.5

tuning the operation predictor and switching the two. In this way, we can guarantee the whole model to converge in a stable process (Wang et al., 2020).

### 3 Experiments

We now detail a series of experiments that we conducted to validate the effectiveness of our proposed planning-LM approach on step-by-step MWP solution generation. Since our focus is on MWP solution generation with explanations, GSM8K (Cobbe et al., 2021) is a good fit for our purpose. This dataset contains 8.5K high-quality and linguistically diverse MWPs, where each MWP has 2-8 solution steps. See Section C for details on data preprocessing.

#### 3.1 Automated Metrics

We need a variety of different metrics to understand the effectiveness of our planning-LM approach. For the final answer, we use the **solve rate** metric to evaluate whether the model generates the final correct answer to each MWP. Since generating meaningful steps is also key, we use the **BLEU** metric (Papineni et al., 2002) to evaluate language generation quality. For intermediate steps, we use the equation match accuracy (**ACC-eq**) metric to evaluate whether a generated step contains a math expression (including numbers) that matches the ground truth. Since LMs generate math equations as strings, we decompose the equation string into

tokens and calculate the token level match rate instead of the overall string match. We also use the operation match accuracy (**ACC-op**) metric to evaluate whether a generated step’s operation label matches the ground truth.

#### 3.2 Human Evaluation

Our proposed planning-LM framework cannot be accurately evaluated using only automated metrics since text similarity metrics such as BLEU do not accurately reflect the mathematical validity of intermediate solution steps. To address this limitation, we implemented a human evaluation protocol with three metrics: *reasoning strategy*, *clear explanation*, and *overall preference*. Ten raters with a good understanding of fundamental mathematics concepts evaluated 50 randomly selected MWPs using the protocol, where their task is to compare two different step-by-step solutions. Each MWP receiving at least three ratings. The full evaluation template can be found in Section G.

#### 3.3 Experimental Settings

We conduct two experiments to verify the effectiveness of our planning-LM framework. In the first, single-step experiment, we input the question and ground-truth solution steps to the model and let it generate the next step and calculate the ACC-eq and ACC-op metrics for each generated step. Since some of the steps are too short, yielding a high variance in BLEU scores, we concatenate all

generated steps and calculate the overall BLEU metric between the ground truth solution and this true history-informed solution. In the second, all-step experiment, we only provide the model with the MWP and ask it to generate all solution steps. We then calculate the solve rate metric to evaluate whether the final answer is correct. We choose GPT-2 (117M parameters) and GPT-2-medium (345M) as our base models and compare the generation results between LM fine-tuning and planning-LM. Meanwhile, we perform another experiment using the ground truth operation prompt as input for planning-LM to generate the next step. The result, an upper bound on the performance of planning-LM, reflects the effectiveness of low-level token-by-token generation in each step, while ACC-eq and ACC-op reflect the effectiveness of high-level mathematical operation planning across steps.

We also conduct the above experiments on encoder-decoder LMs: T5-base(220M) and T5-large(770M). The decoder architecture is the same as GPT-2 models, but instead of treating the question as history input, T5 contains an extra encoder to encode the question and uses cross attention to the question to generate results.

To fairly compare planning-LM with other works on LLM prompting such as chain-of-thought, instead of prompt-tuning on a relatively small LM, we adapt our approach for in-context learning. We select five examples with a specific format  $(Q, P, O_1, S_1, P, O_2, S_2, \dots)$ , i.e., the question followed by a number of prompt-operation-solution triples. We use the examples with GPT-3 (“text-davinci-003”) for in-context learning. An example of the prompt we use is shown in Table 6.

### 3.4 Quantitative Results

#### 3.4.1 Prompt-tuning

Table 1 shows the experimental results for all prompt-tuning based approaches across the two experiments. We see that planning-GPT-2 and planning-T5 with our operation classifier outperform chain-of-thought-tuning on both GPT-2 and T5. We also observe that a similar trend holds for the larger models, GPT-2-medium and T5-large. We highlight that with the planning component, which introduces only around 10K new parameters for the MWP solving task, a base GPT-2 model with 117M parameters performs similarly to a much larger base GPT-2-medium model with 345M parameters. This observation shows that our

Table 3: Solving Rate (%) of different prompting methods on GSM8K. All the in-context prompting methods use text-davinci-003 with 4 examples.

Model	Solve Rate(%)
Standard prompting	21.4
Chain-of-thought	68.5
Planning-LM	72.3

planning approach is highly parameter-efficient for MWP solving. The other observation is that our approach seems to adapt better to decoder-only LMs than to encoder-decoder LMs, even ones with more parameters; T5-base yields almost the same performance as GPT-2, with twice as many parameters.

To validate the effectiveness of each component in our planning-LM approach, we conduct an ablation study on four different components: using prefix or infix prompts, fixed or fine-tuned mathematical operation prompts, instruction prompt mining, and the operation predictor. We see that using infix, fine-tuned mathematical prompts, and the operation predictor improve performance the most across different settings. We also see that infix prompts are significantly better than prefix prompts, which is different from the observation made in prior work (Li and Liang, 2021). One possible explanation is the incompatibility between prefix prompting and step-by-step generation: prefix prompts put the most important instruction at the front of the LM input, making all generated tokens attend to it, which leads to higher operation prediction accuracy but worse generation performance on other tokens.

#### 3.4.2 In-context learning

We conduct experiments by giving in-context prompting examples to GPT-3 in different formats and the result is shown in Table 3. We see that planning-LM yields the best solving rate, significantly higher than other approaches. We further analyze the human evaluation results in Section 3.5.

### 3.5 Human Evaluation Results

Figure 2 shows the distributions of participants’ selections on human evaluation metrics for the generated solutions. We see that solutions generated by planning-LM are significantly better than those produced by chain-of-thought on all three metrics, proving that our approach leads to solutions with more precise language and better problem solving

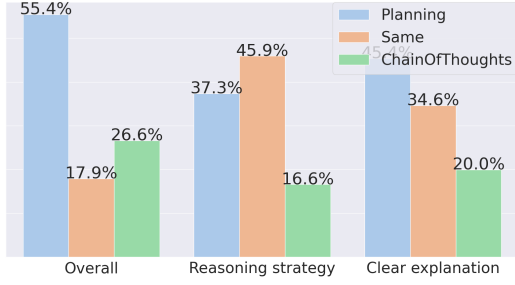


Figure 2: The distributions of participants’ selections in our human evaluation experiment.

strategies. Providing math operation information to the LM as hints on the next step also help the model to generate more clear and sound explanations in the intermediate solution steps.

### 3.6 Qualitative Analysis

Table 4 shows two examples that compare the full step-by-step solutions generated by our planning-LM approach and chain-of-thought prompting. For Example 1, we see that although chain-of-thought happens to produce the correct answer, the reasoning starts to fall apart at Step 3. It generated the correct final answer only because the question mentioned rounding the answer to the nearest integer; however, its intermediate answer 1.33 is wrong. For Example 2, the answer generated by the chain-of-thought does not have detailed wording explanations, whereas planning LM’s solution has details of each step of the solving strategy, making the solution much easier to understand.

Perhaps surprisingly, we observe that planing-LM can generate multiple solutions if it predicts a different math operation in the next step compared to the ground truth solution. Therefore, we conduct a follow-up experiment by giving the model a hand-crafted plan via operation prompts to see whether it can generate an alternative correct solution strategy.

Table 5 further demonstrates that our approach can generate multiple correct solution paths for the same problem. For example, feeding Plans I and II enables the model to generate the correct final answer among the four strategies we used; the generated solutions follow the operation steps given, indicating that the model has some reasoning ability and can extract some meaningful patterns from data. Plan III results in a flawed solution and Plan IV failed since we do not have an operation class that matched the step. For plan III, the first step,  $[n + n + \dots]$ , is not seen often enough in the training data. For plan IV,  $(n + n) \times n$  is not seen in

the training data either. However, we note that in this case, using the closest operation,  $[n + n \times n]$ , results in a solution that gets very close to the correct final answer. These results suggest that a better representation of the operation prompt is crucial for future work since our current approach is limited to a finite number of predefined operations; a prompt operation *generator* rather than classifier could be a better choice for a wide variety of mathematical operations. We also note that this flexibility gives our planning-LM approach potential to be useful in real-world applications. For example, these solution plan controls may encourage math students to explore different solution strategies and be more creative.

## 4 Related work

**MWP solver** A large body of recently proposed MWP solvers parses an MWP into its underlying equation, which has been a very active research area with a plethora of related work. These works differ mainly in the technical approaches which broadly fall in three categories. First, some works explore MWP solving via reinforcement learning, which rewards the model with the correct answer generated (Huang et al., 2018; Wang et al., 2018b). RL methods generally requires a sizable dataset and can be unstable to train, which may not be suitable for most MWP datasets that are only of modest sizes. Second, some works exploit the combination of symbolic- and neural-network-based approaches, e.g., by combining a pre-defined symbolic patterns such as solution templates (Wang et al., 2019) and symbolic tree structures of equations (Xie and Sun, 2019; Li et al., 2020; Qin et al., 2020; Wang et al., 2018a; Wu et al., 2020; Zhang et al., 2021). These methods can be significantly constrained by these patterns and it may be challenging to generalize them to other MWPs whose solutions are not expressed by these patterns.

Lastly, some works build on large LMs (LLMs) via special fine-tuning or inference techniques. Chain-of-thought prompting (Wei et al., 2022) prompts LLMs to generate intermediates steps before reaching the final answer. Cobbe et al. (2021) fine-tunes a model as a verifier and applies the verifier to rank outputs in the decoding phase. Wang et al. (2022) are using a majority vote among outputs to select the best answer. Lewkowycz et al. (2022) fine-tunes an LLM by a large collection of math-specific datasets combining existing tech-

Table 4: Demonstrations of generated solutions comparing planning-LM and chain-of-thought. Question 1 shows the intermediate step of chain-of-thought has wrong reasoning but still reaches the final answer. Question 2 shows that planning-LM results in a better reasoning strategy since the calculation process is simple and more concrete.

**1.Question:** Antoine’s french onion soup recipe calls for 2 pounds of onions. He likes to double that amount. His soup serves 6 people. The onions are currently on sale for \$2.00 a pound. He also needs 2 boxes of beef stock, that are also on sale for \$2.00 a box. What is the cost per serving? (Round to the nearest integer.)

**Chain-of-thought:** The cost of the onions is 2 pounds times \$2.00 per pound, which is  $2 \times 2 = 4$  dollars. The cost of the beef stock is 2 boxes times \$2.00 per box, which is  $2 \times 2 = 4$  dollars. The total cost of the ingredients is  $4 + 4 = 8$  dollars. The cost per serving is 8 divided by 6 servings, which is  $8/6 = 1.33$  dollars.

Rounded to the nearest integer, the cost per serving is \$2.  
**PlanningLM:** Antoine needs 4 pounds of onions, which costs  $4 \times 2 = 8$  dollars. He also needs 2 boxes of beef stock, which costs  $2 \times 2 = 4$  dollars. The total cost for the ingredients is  $8 + 4 = 12$  dollars. The cost per serving is  $12/6 = 2$  dollars.

**2.Question:** Judy teaches 5 dance classes, every day, on the weekdays and 8 classes on Saturday. If each class has 15 students and she charges \$15.00 per student, how much money does she make in 1 week?

**Chain-of-thought:** On the weekdays, Judy teaches 5 classes per day, with 15 students in each class. This means she has  $5 \times 15 = 75$  students per day. At a rate of \$15 per student, she makes  $75 \times 15 = 1125$  dollars per day. For the weekdays, she makes  $5 \times 1125 = 5625$  dollars. On Saturday, Judy teaches 8 classes with 15 students in each class. This means she has  $8 \times 15 = 120$  students. At a rate of \$15 per student, she makes  $120 \times 15 = 1800$  dollars. For the week, Judy makes  $5625 + 1800 = 7425$  dollars. Rounded to the nearest integer, the cost per serving is \$2.

**PlanningLM:** Judy teaches 5 classes per weekday, so  $5 \times 5 = 25$  classes. Judy teaches 8 classes on Saturday, so  $8 \times 1 = 8$  classes. The total number of classes she teaches in a week is  $25 + 8 = 33$  classes. Each class has 15 students, so  $15 \times 33 = 495$  students. Judy charges \$15.00 per student, so  $495 \times 15 = \$7,425$ .

niques. There are also some extension works based on CoT, like the least-to-most prompting (Zhou et al., 2022) that decomposes the complicated question into a single-hop question; STaR (Zelikman et al., 2022) iterative rationale generation using a small number of examples and a large dataset. Our work differs from previous studies by not only prioritizing the final solution accuracy but also emphasizing the generation quality of individual solution steps. Additionally, we introduce a novel hierarchical planning method for fine-tuning, in contrast to previous approaches that rely solely on ordinary language modeling techniques.

Table 5: Qualitative examples of using our planning-LM to plan for different but valid solution strategies to achieve the same correct result for a given MWP. Plan IV failed since we do not have an exactly operation class that matched the step.

**Question:** Two trains leave San Rafael at the same time. They begin traveling westward, both traveling for 80 miles. The next day, they travel northwards, covering 150 miles. What’s the distance covered totally in the two days? Ans: 460

**Plan I:**  $[n + n] \rightarrow [n * n] \rightarrow [Ans]$

**Result:** The total distance covered in the two days is  $80 + 150 = 230$  miles. The total distance covered in the two days is  $230 * 2 = 460$ . Answer is 460.

**Plan II:**  $[n * n] \rightarrow [n * n] \rightarrow [n + n] \rightarrow [Ans]$

**Result:** The first train covers  $80 * 2 = 160$  miles. On the second day, the trains covered 2 trains  $* 150$  miles/train = 300 miles. Total distance covered is  $160 + 300 = 460$  miles. Answer is 460.

**Plan III:**  $[n + n + \dots] \rightarrow n + n \rightarrow [Ans]$

**Result:** The total distance covered in the two days is  $80 + 80 + 150 = 310$  miles. The total distance covered in one day is  $310 + 150 = 460$  miles. Answer is 460.

**Plan IV:**  $[n + n * n] \rightarrow [Ans]$

**Result:** The total distance covered by trains in the two days is  $150 + 80 * 2 = 310$  miles. Answer is 310.

**Controllable text generation** Given the rise of LLMs, controllable generation methods that guide these large models to generate desirable content and avoid potential pitfalls such as bias (Bender et al., 2021) has been a recent research trend. These controllable generation methods generally fall into two categories. Works in the first category modify the token distribution at each time step to achieve controllable generation via gradient-based methods (Dathathri et al., 2020), external classifier-based methods (Krause et al., 2020; Liu et al., 2021a), or resampling (Clark et al., 2020; Bhattacharyya et al., 2021; Bakhtin et al., 2021). Works in the second category fine-tunes the LLM via either language modeling (Shirish Keskar et al., 2019) or reinforcement learning (Khalifa et al., 2021). These works focus on controllable generation for natural language and study nuanced control attributes such as topic and sentiment that can only be matched implicitly. In contrast, our work focuses differently on both natural and mathematical language, which involves control attributes, e.g., math operation hints in the form of equations that need to be matched exactly.



## 5 Conclusions

In this paper, we addressed the new problem of performing fine-grained, step-by-step controllable solution generation for math word problems. We proposed an approach combining planning and language models to generate interpretable solution steps. Our approach leverages pre-trained language models in two ways: at each step, plan the mathematical operation to be applied, followed by using these plans as prompts to control the token-by-token generation of each step. We demonstrated that with minimal additional parameters introduced, our approach significantly improves math word problem-solving performance over simply fine-tuning language models. We also showed that due to the interpretability and high controllability of operation prompts, we can use our approach to generate solutions with alternative strategies by giving it different solution plans. Future work can further explore generating an entire solution path by predicting math operators for each step and revising the plan after each step is generated. We can also explore the application of our approach in real-world educational settings, e.g., for open-ended answer scoring (Lan et al., 2015; Zhang et al., 2022).

## 6 Limitations

First, our work applies hand-crafted action labels as operation hints, which leads to some limitations to represent more complex operation steps. For the future work, we can use a generator instead of a classifier to generate a more flexible set of operation prompts, making them more representative and meaningful. Secondly, due to the high controllable generation of our approach, if our approach yields a wrong operation step prediction, it would further mislead the intermediate step generation. To eliminate the drawback where inaccurately generated operation prompts would mislead the next step, we can apply a verifier (Cobbe et al., 2021) to evaluate the reliability of the generated operation prompts. When the reliability is low, we ditch the operation prompt to prevent it from guiding the model into an incorrect path.

## 7 Ethics Statement

Currently, most existing works leverage the capability of generating intermediate reasoning steps of large, pre-trained language models for either understanding the model’s behaviors (e.g., models’

moral judgments (Jin et al., 2022)) or improving their problem-solving accuracies (e.g., MWP solving (Lewkowycz et al., 2022)). Few works focus on the quality of the generated intermediate reasoning steps themselves. These generated steps have potentially significant real-world applications, such as providing feedback automatically in large-scale education scenarios, but they are not yet of high enough quality to be readily utilized in practice. Our work contributes to the important direction in making such generated intermediate steps more accurate, coherent, and high-quality. However, language models equipped with our approach may still generate intermediate steps that are unreasonable, even though it improves upon existing approaches. These unreasonable generated steps may be misleading to students when they are learning, posing a potential risk to their usage. As a result, more work is required before our approach can be readily deployed in practice. We believe that, in its current form, our work is best suitable for use with experts, i.e., education subject matter experts or instructors to help them write solution steps for new MWPs in a more efficient manner.

## Acknowledgement

The authors thank the NSF (under grants 1917713, 2118706, 2202506, 2215193, 2237676) for partially supporting this work.

## References

- Anton Bakhtin, Yuntian Deng, Sam Gross, Myle Ott, Marc’Aurelio Ranzato, and Arthur Szlam. 2021. Residual energy-based models for text. *J. Mach. Learn. Res.*, 22(40):1–41.
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proc. ACM Conf. Fairness Accountability Transparency*, page 610–623.
- Sumanta Bhattacharyya, Amirmohammad Rooshenas, Subhajit Naskar, Simeng Sun, Mohit Iyyer, and Andrew McCallum. 2021. Energy-based reranking: Improving neural machine translation using energy-based models. In *Proc. Annu. Meeting Assoc. Comput. Linguistics and Int. Joint Conf. Natural Lang. Process.*, pages 4528–4537.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie

- Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2021. [On the opportunities and risks of foundation models](#).
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Yixuan Cao, Feng Hong, Hongwei Li, and Ping Luo. 2021. A bottom-up dag structure extraction model for math word problems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 39–46.
- Xinyun Chen, Chen Liang, Adams Wei Yu, Denny Zhou, Dawn Song, and Quoc V. Le. 2020. [Neural symbolic reader: Scalable integration of distributed and symbolic representations for reading comprehension](#). In *International Conference on Learning Representations*.
- Chowdhery, Aakanksha and others. 2022. [Palm: Scaling language modeling with pathways](#). *arXiv preprint <https://arxiv.org/abs/2204.02311>*.
- Kevin Clark, Minh-Thang Luong, Quoc Le, and Christopher D. Manning. 2020. Pre-training transformers as energy-based cloze models. In *Proc. Conf. Empirical Methods Natural Lang. Process.*, pages 285–294.
- 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*.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. Plug and play language models: A simple approach to controlled text generation. In *Proc. Int. Conf. Learn. Representations*.
- Edward A Feigenbaum and Julian Feldman, editors. 1995. *Computers and Thought*. MIT Press, London, England.
- Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. 2021. [WARP: word-level adversarial reprogramming](#). *CoRR*, abs/2101.00121.
- Danqing Huang, Jing Liu, Chin-Yew Lin, and Jian Yin. 2018. Neural math word problem solver with reinforcement learning. In *Proc. ACL*, pages 213–223.
- Zhanming Jie, Jierui Li, and Wei Lu. 2022. [Learning to reason deductively: Math word problem solving as complex relation extraction](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5944–5955, Dublin, Ireland. Association for Computational Linguistics.
- Zhijing Jin, Sydney Levine, Fernando Gonzalez, Ojasv Kamal, Maarten Sap, Mrinmaya Sachan, Rada Mihalcea, Josh Tenenbaum, and Bernhard Schölkopf. 2022. [When to make exceptions: Exploring language models as accounts of human moral judgment](#).
- Jeffery D. Karpicke. 2012. Retrieval-based learning: Active retrieval promotes meaningful learning. *Current Directions Psychol. Sci.*, 21(3):157–163.
- Muhammad Khalifa, Hady Elsahar, and Marc Dymetman. 2021. A distributional approach to controlled text generation. In *Proc. Int. Conf. Learn. Representations*.
- Kenneth R. Koedinger, Jihee Kim, Julianna Zhuxin Jia, Elizabeth A. McLaughlin, and Norman L. Bier. 2015. Learning is not a spectator sport: Doing is better than watching for learning from a mooc. In *Proc. Conf. Learn. Scale*, pages 111–120.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*.
- Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. 2020. GeDi: Generative Discriminator Guided Sequence Generation. *arXiv e-prints*.
- Vivek Kumar, Rishabh Maheshwary, and Vikram Pudi. 2022. [Practice makes a solver perfect: Data augmentation for math word problem solvers](#). In *Proceedings*

- of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4194–4206, Seattle, United States. Association for Computational Linguistics.
- Andrew S Lan, Divyanshu Vats, Andrew E Waters, and Richard G Baraniuk. 2015. Mathematical language processing: Automatic grading and feedback for open response mathematical questions. In *Proceedings of the ACM conference on learning@scale*, pages 167–176.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). *CoRR*, abs/2104.08691.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay 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](#).
- Shucheng Li, Lingfei Wu, Shiwei Feng, Fangli Xu, Fengyuan Xu, and Sheng Zhong. 2020. Graph-to-tree neural networks for learning structured input-output translation with applications to semantic parsing and math word problem. In *Proc. EMNLP*, pages 2841–2852.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). *CoRR*, abs/2101.00190.
- Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2022. On the advance of making language models better reasoners. *arXiv preprint arXiv:2206.02336*.
- Zhongli Li, Wenxuan Zhang, Chao Yan, Qingyu Zhou, Chao Li, Hongzhi Liu, and Yunbo Cao. 2021. [Seeking patterns, not just memorizing procedures: Contrastive learning for solving math word problems](#). *CoRR*, abs/2110.08464.
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021a. DExperts: Decoding-time controlled text generation with experts and anti-experts. In *Proc. Annu. Meeting Assoc. Comput. Linguistics and Int. Joint Conf. Natural Lang. Process.*, pages 6691–6706.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021b. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing](#). *CoRR*, abs/2107.13586.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021c. [GPT understands, too](#). *CoRR*, abs/2103.10385.
- Aman Madaan and Amir Yazdanbakhsh. 2022. [Text and patterns: For effective chain of thought, it takes two to tango](#).
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: A method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL ’02, page 311–318, USA. Association for Computational Linguistics.
- Jinghui Qin, Lihui Lin, Xiaodan Liang, Rumin Zhang, and Liang Lin. 2020. Semantically-aligned universal tree-structured solver for math word problems. In *Proc. EMNLP*, pages 3780–3789.
- Zhihong Shao, Fei Huang, and Minlie Huang. 2022. [Chaining simultaneous thoughts for numerical reasoning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2533–2547, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- 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.
- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. CTRL: A Conditional Transformer Language Model for Controllable Generation. *arXiv e-prints*.
- Sowmya S. Sundaram, Sairam Gurajada, Marco Fisichella, Deepak P, and Savitha Sam Abraham. 2022. [Why are NLP models fumbling at elementary math? A survey of deep learning based word problem solvers](#). *CoRR*, abs/2205.15683.
- Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. 2022. [Solving math word problems with process- and outcome-based feedback](#).
- Lieven Verschaffel, Stanislaw Schukajlow, Jon Star, and Wim Van Dooren. 2020. Word problems in mathematics education: a survey. *ZDM*, 52(1):1–16.
- Candace A. Walkington. 2013. Using adaptive learning technologies to personalize instruction to student interests: The impact of relevant contexts on performance and learning outcomes. *J. Educ. Psychol.*, 105(4):932–945.
- Jianhong Wang, Yuan Zhang, Tae-Kyun Kim, and Yun-jie Gu. 2020. [Modelling hierarchical structure between dialogue policy and natural language generator with option framework for task-oriented dialogue system](#). *CoRR*, abs/2006.06814.
- Lei Wang, Yan Wang, Deng Cai, Dongxiang Zhang, and Xiaojiang Liu. 2018a. Translating a math word problem to a expression tree. In *Proc. EMNLP*, pages 1064–1069.

- Lei Wang, Dongxiang Zhang, Lianli Gao, Jingkuan Song, Long Guo, and Heng Tao Shen. 2018b. Mathdqn: Solving arithmetic word problems via deep reinforcement learning. In *Proc. AAAI*, pages 5545–5552.
- Lei Wang, Dongxiang Zhang, Jipeng Zhang, Xing Xu, Lianli Gao, Bing Tian Dai, and Heng Tao Shen. 2019. Template-based math word problem solvers with recursive neural networks. In *Proc. AAAI*, volume 33, pages 7144–7151.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. [Self-consistency improves chain of thought reasoning in language models](#).
- Yan Wang, Xiaojiang Liu, and Shuming Shi. 2017. Deep neural solver for math word problems. In *Proc. EMNLP*, pages 845–854.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. 2022. [Chain of thought prompting elicits reasoning in large language models](#). *CoRR*, abs/2201.11903.
- Qinzhao Wu, Qi Zhang, Jinlan Fu, and Xuanjing Huang. 2020. A knowledge-aware sequence-to-tree network for math word problem solving. In *Proc. EMNLP*, pages 7137–7146.
- Zhipeng Xie and Shichao Sun. 2019. A goal-driven tree-structured neural model for math word problems. In *Proc. IJCAI*.
- Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. [Bartscore: Evaluating generated text as text generation](#). *CoRR*, abs/2106.11520.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488.
- Mengxue Zhang, Sami Baral, Neil Heffernan, and Andrew Lan. 2022. Automatic short math answer grading via in-context meta-learning. *arXiv preprint arXiv:2205.15219*.
- Qiyuan Zhang, Lei Wang, Sicheng Yu, Shuohang Wang, Yang Wang, Jing Jiang, and Ee-Peng Lim. 2021. [NOAHQA: Numerical reasoning with interpretable graph question answering dataset](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4147–4161, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021. [Factual probing is \[MASK\]: learning vs. learning to recall](#). *CoRR*, abs/2104.05240.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Olivier Bousquet, Quoc Le, and Ed Chi. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*.
- Yanyan Zou and Wei Lu. 2019. [Text2math: End-to-end parsing text into math expressions](#).



## A Hyper-parameters

We use a learning rate of  $5e-5$ , a batch size of 8, and 10 epochs for all training processes. We set “what is the next operation?” as our instruction prompt and apply *calculators* to avoid calculation errors and *greedy decoding* during token generation. Model training is carried out on an NVIDIA RTX 3090 GPU.

## B Setting for generation

In order to have consistent results and fair comparison, we apply a greedy-decoding procedure with  $\text{temp}=0$  for all of the generation process during testing.

## C Data preprocessing

We detail how to obtain our step operation here. For each MWP, we split the solution into steps according to the period symbol “.” at the end of sentences. We restrict ourselves to the top-20 most frequent mathematical operations after merging some operations that have a similar meaning, e.g.,  $[n + n + n]$  and  $[n + n + n + n]$  are both labeled as “multi-step addition” to avoid highly infrequent operations. Check table 7 for full descriptions of operation actions. We applied packages ‘NLTK’ and ‘Spacy’ for data preprocessing.

## D More Details of prompt tuning design

### D.1 Operation prompts

We initialize the embedding of each math operation token as the original pre-trained LM’s embedding of the mathematical operator token instead of initializing them randomly (Liu et al., 2021c). For example, we initialize the operations action token  $[n + n]$  with the same value as embedding of the “+” token in the pre-trained model. For operation classes that contain multiple operations, we initialize the embedding to the mean of all operation embeddings involved. We do this since initializing a new token with related embeddings has been proven to be effective on speeding up the training process of LM-based models (Li and Liang, 2021; Zhong et al., 2021; Lester et al., 2021; Hambardzumyan et al., 2021; Liu et al., 2021b).

### D.2 Prompt mining through paraphrasing

For the instruction prompt, finding good prompts is an art that takes time and experience (Liu et al., 2021b). Thus, we apply prompt mining through paraphrasing by first starting with a seed prompt (e.g. “The next step operation is: ”) and paraphrase it into a set of other candidate prompts with similar meaning (Yuan et al., 2021). Then, we tune the model with these candidates by treating them as hyper-parameters and select the one that performs best on the target task. We find that anchor tokens (e.g. “?”) are helpful and leads to good performance, which is consistent with prior work (Liu et al., 2021c).

## E Architecture for text-to-text language modeling

See figure 3

## F In-context Prompting design

Table 6: The complete inputs of in-context prompting for planning-LM

---

QUESTION: The pet shop grooms dogs. It takes 30 minutes to groom a poodle. It takes half as much time to groom a terrier as it takes to groom a poodle. They do not groom cats. If the pet shop grooms 3 poodles and 8 terriers, what is the total length of time it will take in minutes?
SOLUTION:
: [*=]
At a rate of 30 minutes per poodle, 3 poodles will take $3*30=\ll 3*30=90 \gg 90$ minutes.
: [/=]
If it takes half as much time to groom a terrier as it takes to groom a poodle, it takes $30/2=\ll 30/2=15 \gg 15$ minutes per terrier.
: [*=]

---

At a rate of 15 minutes per terrier, 8 terriers will take  $\ll 15 \times 8 = 120 \gg 120$  minutes.

: [+=]

Thus, in total, it will take the pet shop  $90 + 120 = \ll 90 + 120 = 210 \gg 210$  minutes to groom 3 poodles and 8 terriers.

: [end]

boxed{210}

QUESTION: Darnell has 1000 square feet of fabric that he's using to make mini flags. He makes square flags that are 4 feet by 4 feet, wide rectangular flags that are 5 feet by 3 feet, and tall rectangular flags that are 3 feet by 5 feet. He has already made 16 square flags, 20 wide flags, and 10 tall flags. How many square feet of fabric does he have left?

SOLUTION:

: [\*=]

The square flags use 16 square feet of fabric per flag because four times four equals  $\ll 4 \times 4 = 16 \gg 16$

: [\*=]

The wide flags each use 15 feet of fabric per flag because five times three equals  $\ll 5 \times 3 = 15 \gg 15$

: [\*=]

The tall flags each use 15 feet of fabric per flag because five times three equals  $\ll 5 \times 3 = 15 \gg 15$

: [\*=]

He has used 256 square feet for the square flags because 16 times 16 equals  $\ll 16 \times 16 = 256 \gg 256$

: [\*=]

He has used 300 square feet for the wide flags because 20 times 15 equals  $\ll 20 \times 15 = 300 \gg 300$

: [\*=]

He has used 150 square feet for the tall flags because ten times 15 equals  $\ll 15 \times 10 = 150 \gg 150$ .

: [++=]

He has used 706 square feet because 256 plus 300 plus 150 equals  $\ll 256 + 300 + 150 = 706 \gg 706$ .

: [-=]

He has 294 square feet left because 1,000 minus 706 equals  $\ll 1000 - 706 = 294 \gg 294$

: [end]

boxed{294}

QUESTION: Mr. Smith takes his wife, his parents, and his 3 children to a seafood buffet. The price for the adult buffet is \$30. The price for the children's buffet is \$15. Senior citizens get a 10% discount. How much does Mr. Smith spend on the buffet for his entire family?

SOLUTION:

: [\*=]

The buffet cost for Mr. Smith and his wife is  $\$30 \times 2 = \$\ll 30 \times 2 = 60 \gg 60$ .

: [\*=]

The buffet cost for his children is  $\$15 \times 3 = \$\ll 15 \times 3 = 45 \gg 45$ .

: [\*\*.]=]

The buffet cost for 1 senior citizen is  $\$30 \times 90\% = \$\ll 30 \times 90 \times .01 = 27 \gg 27$ .

: [\*=]

Buffet cost for the grandparents is  $\$27 \times 2 = \$\ll 27 \times 2 = 54 \gg 54$ .

: [++=]

Buffet cost for the entire family is  $\$60 + \$45 + \$54 = \$\ll 60 + 45 + 54 = 159 \gg 159$

: [end]

boxed{159}

QUESTION: Jenny's local library received 50 visitors on Monday. On Tuesday, the library received twice that number. On the remaining days of the week, an average of 20 visitors went to the library. What's the total number of visitors who went to the library that week?

SOLUTION:

: [\*=]

On Monday, there were 50 visitors. On Tuesday, there were twice as many, so  $2 \times 50 = \ll 2 \times 50 = 100 \gg 100$  visitors

: [+=]

The total number of visitors after Tuesday will be  $100 + 50 = \ll 100 + 50 = 150 \gg 150$  visitors.

: [\*=]

For the remaining five days, an average of 20 visitors attended, giving a total of  $5 \times 20 = \ll 5 \times 20 = 100 \gg 100$  visitors.

: [+]=

The total number of visitors who visited the library for the week was  $100+150 = \langle 100+150=250 \rangle 250$  visitors.

: [end]

boxed{250}

QUESTION: James decides to build a tin house by collecting 500 tins in a week. On the first day, he collects 50 tins. On the second day, he manages to collect 3 times that number. On the third day, he collects 50 tins fewer than the number he collected on the second day. If he collects an equal number of tins on the remaining days of the week, what's the number of tins he collected each day for the rest of the week?

SOLUTION:

: [\*]=

On the second day, he collected 3 times the number of tins he collected on the first day, which is  $3*50 = \langle 3*50=150 \rangle 150$  tins.

: [-]=

On the third day, he collected 50 tins fewer than the second day, which is  $150-50 = \langle 150-50=100 \rangle 100$  tins

: [++]=

The total for the three days is  $150+100+50 = \langle 150+100+50=300 \rangle 300$  tins.

: [-]=

To reach his goal, he still needs  $500-300 = \langle 500-300=200 \rangle 200$  tins.

: [/]=

Since the total number of days left in the week is 4, he'll need to collect  $200/4 = \langle 200/4=50 \rangle 50$  tins per day to reach his goal

: [end]

boxed{50}

QUESTION: Lilah's family gallery has 400 photos. On a two day trip to the Grand Canyon, they took half as many photos they have in the family's gallery on the first day and 120 more photos than they took on the first day on the second day. If they added all these photos to the family gallery, calculate the total number of photos in the gallery.

SOLUTION:

: [/]=

On their first day at the grand canyon, the family took half as many photos as the ones they have in the gallery, meaning they took  $1/2*400 = \langle 400/2=200 \rangle 200$  photos.

: [+]=

The total number of photos, if they add the ones they took on the first day to the family's gallery, is  $400+200 = \langle 400+200=600 \rangle 600$

: [+]=

On the second day, they took 120 more photos than they took on the first day, a total of  $200+120 = \langle 200+120=320 \rangle 320$  photos.

: [+]=

After adding the photos they took on the second day to the galley, the number of photos will be  $600+320 = \langle 600+320=920 \rangle 920$

: [end]

boxed{920}

## G Human Evaluation template

Figure 4 shows the template for human evaluation. I collected data through google forms platforms. Chain-Of-Thoughts or Planing-LM solutions would randomly assign solution A or solution B so the participants cannot identify which one is which one. See figure 3

## H List of all hand-crafted operations classes

Details are in table 7

## I Examples of control generation

Table 8 shows the generated step apply different operation prompts on same input. This table demonstrates the generated results from applying different operation prompts with the same input to the model. We

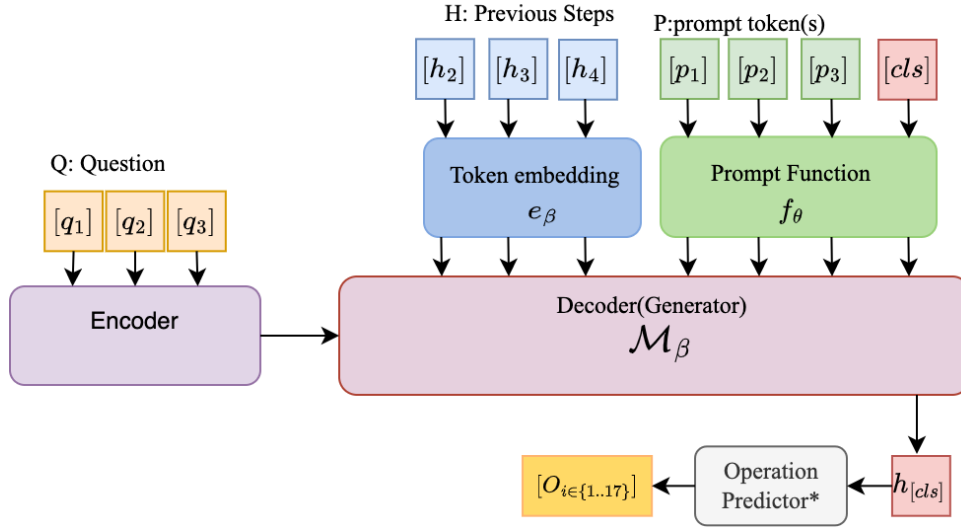


Figure 3: An overview architecture for text-to-text transformer based language model such as T5. The question is not considered as history; instead will feed into an encoder and also have cross-attention to the generated steps tokens.

Table 7: Appendix C: list of all hand-crafted operations classes

no.	shortcut	description	no.	shortcut	description
1	[n+n]	one-step addition	12	[n-n*n]	multiplication then subtraction
2	[n-n]	one-step subtraction	11	[n*(n/n)]	multiplication by a fraction
3	[n*n]	one-step multiplication	13	[(n/n)-(n/n)]	fraction subtraction
4	[n/n]	one step division	14	[(n/n)+(n/n)]	fraction addition
5	[n+n+. . .]	multi step addition	11	[(n/n)*(n/n)]	fraction multiplication
6	[n-n-. . .]	multi-step subtraction	16	[mixed]	other combination
7	[n*n*. . .]	multi-step multiplication	17	[ans]	solution found, end the whole generation
8	[n+n*n]	multiplication then addition	18	[statement]	involve no math calculation and only textual explanation
9	[n+n-n]	addition then subtraction	19	[assign]	assign a value to a paramter
10	[n+n/n]	division then addition	20	[define]	define a parameter

observe that when the operation prompt is logical and aligned with solving the question, the generated result follows the guidance given by the operation prompt wit higher probability. In contrast, when the operation prompt does not make sense, there is the chance that the generated result will not obey its directions. Details in table 8



Table 8: Generated step apply different operation prompts on same input.

Input: Maggie spent a quarter of her money, while Riza spent one-third of her money.  
They each had \$60. How much money do the two of them have left?

no.	action	generated step
1	[n+n]	Riza and Maggie had $\$60 + \$60 = \$120$ ✓
2	[n-n]	Riza had $\$60 - \$60 = \$0$ left after spending some money on Maggie. ✓
3	[n*n]	Riza spent $\$60 \times 0.33 = \$20$ ✓
4	[n/n]	Maggie had $\$60/4 = \$15$ left. ✓
5	[n+n+. . .]	Riza spent $\$60 + \$60 + \$60 = \$180$ . ✓
6	[n-n-. . .]	Riza had $\$60 - \$60 - \$30 = -\$30$ ✓
7	[n*n*. . .]	Riza spent $\$60 \times 1/3 = \$60$ .
8	[n+n*n]	Maggie had $\$60 + \$60 \times 0.25 = \$56$ left after spending a quarter of her money. ✓
9	[n+n-n]	Maggie had $\$60 - \$60 = \$0$ left after spending a quarter of her money.
10	[n+n/n]	Riza had $\$60 + \$60/3 = \$80$ left. ✓
11	[n-n*n]	Riza had $\$60 - \$30 = \$30$ left after spending a quarter of her money.
12	[n*(n/n)]	Riza spent $\$60 \times 1/3 = \$20$ . ✓
13	[(n/n)-(n/n)]	Riza and Maggie had $\$60/4 - \$60/3 = \$5$ left. ✓
14	[(n/n)+(n/n)]	Riza and Maggie total had $\$60/4 + \$60/3 = \$35$ left. ✓
11	[(n/n)*(n/n)]	Riza had $\$60 * 1/3 = \$20$ left.
16	[mixed]	Riza had $\$60/4 - \$60/3 = \$15$ left.
17	[ans]	The answer is \$60 ✓
18	[statement]	Riza had \$60 in total. ✓
19	[assign]	Riza had $x = \$60$ . ✓
20	[define]	Let x be the money Riza had ✓

## Math evaluation 1-10

1. Read the math question and corresponding solutions.
2. Rate which one is better based on four different criteria.

Switch account



Question: 273

Terri is knitting a sweater with two sleeves, a collar, and a decorative rosette. The body of the sweater takes 900 stitches to complete, the collar takes a tenth of that number of stitches, and the rosette takes twice as many as the collar. The whole sweater is an 1800-stitch project. How many stitches does each sleeve take?

A:

The body of the sweater takes 900 stitches, the collar takes a tenth of that number, which is  $900/10 = 90$  stitches, and the rosette takes twice as many as the collar, which is  $90 \times 2 = 180$  stitches. The total number of stitches for the sweater is  $900 + 90 + 180 = 1170$  stitches. The remaining stitches are  $1800 - 1170 = 630$  stitches, which will be divided equally between the two sleeves, so each sleeve will take  $630/2 = 315$  stitches.

B:

The body of the sweater takes 900 stitches to complete. The collar takes a tenth of that number of stitches, which is  $900/10 = 90$  stitches. The rosette takes twice as many as the collar, which is  $90 \times 2 = 180$  stitches. The total number of stitches for the sweater is 1800, so the remaining stitches are  $1800 - 900 - 90 - 180 = 630$  stitches. Each sleeve takes  $630/2 = 315$  stitches.

	A	Same	B
Overall prefer	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Easier to follow	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Better reasoning strategy	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Clearer explanation	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Figure 4: The template for human evaluation.