

A Diversity-Enhanced Knowledge Distillation Model for Practical Math Word Problem Solving

Yi Zhang^{a,b}, Guangyou Zhou^c, Zhiwen Xie^c, Jinjin Ma^a and Jimmy Xiangji Huang^d

^aFaculty of Artificial Intelligence in Education, Central China Normal University, Wuhan, 430079, Hubei, China

^bSchool of Information and Safety Engineering, Zhongnan University of Economics and Law, Wuhan, 430073, Hubei, China

^cHubei Provincial Key Laboratory of Artificial Intelligence and Smart Learning & School of Computer, Central China Normal University, Wuhan, 430079, Hubei, China

^dInformation Retrieval and Knowledge Management Research Lab, York University, Toronto, Canada

ARTICLE INFO

Keywords:

Knowledge distillation

Math word problem

Variational auto-encoder

Question answering

ABSTRACT

Math Word Problem (MWP) solving is a critical task in natural language processing, has garnered significant research interest in recent years. Various recent studies heavily rely on Seq2Seq models and their extensions (e.g., Seq2Tree and Graph2Tree) to generate mathematical equations. While effective, these models struggle to generate diverse but counterpart solution equations, limiting their generalization across various math problem scenarios. In this paper, we introduce a novel Diversity-enhanced Knowledge Distillation (DivKD) model for practical MWP solving. Our approach proposes an adaptive diversity distillation method, in which a student model learns diverse equations by selectively transferring high-quality knowledge from a teacher model. Additionally, we design a diversity prior-enhanced student model to better capture the diversity distribution of equations by incorporating a conditional variational auto-encoder. Extensive experiments on four MWP benchmark datasets demonstrate that our approach achieves higher answer accuracy than strong baselines while maintaining high efficiency for practical applications.

1. Introduction

The ability for mathematical reasoning has long been recognized as a fundamental challenge for computers (Bobrow, 1964). Math word problem (MWP) solving aims to generate solutions from mathematical problems expressed in natural language. Solving MWP requires natural language understanding and mathematical reasoning skills, which have garnered significant attention from the fields of natural language processing (NLP) and smart education.

Researchers developed various methods to address MWP solving tasks, including statistical machine learning methods (Kushman et al., 2014; Hosseini et al., 2014; Mitra and Baral, 2016; Roy and Roth, 2018), semantic parsing methods (Shi et al., 2015; Koncel-Kedziorski et al., 2015; Huang et al., 2017), and deep learning methods (Wang et al., 2017; Xiao et al., 2023; Zhang et al., 2024). Traditional statistical machine learning and semantic parsing methods require manually crafted feature engineering or template design, which are challenging to scale to large datasets. Recently, deep learning models have emerged as a promising paradigm for MWP solving. Leveraging the success of end-to-end models in NLP, researchers introduced Sequence-to-Sequence (Seq2Seq) models for MWP tasks (Wang et al., 2017, 2019; Li et al., 2019). Seq2Seq models utilize an Encoder-to-Decoder (Enc2Dec) framework to translate the input problem description into a mathematical equation. However, these models often overlook the structural information of equations, potentially generating invalid equations that cannot be computed. To better capture equation structure, some studies have enhanced Seq2Seq models by incorporating tree-based decoders, such as Seq2Tree models (e.g., GTS) (Xie and Sun, 2019) and Graph2Tree models (e.g., Graph2Tree-Z (Zhang et al., 2020b)).

Although existing Enc2Dec models (e.g., Seq2Seq, Seq2Tree and Graph2Tree) have achieved remarkable progress in solving MWPs, they are limited in modeling the diversity of solution equations due to data and model limitations. On the one hand, benchmark datasets (e.g., Math23K (Wang et al., 2017) and MAWPS (Koncel-Kedziorski et al., 2016)) typically provide only one ground-truth solution equation for each math problem, even though numerous alternative correct equations can solve the same problem. For instance, in the math problem shown in Figure 1, the ground-truth

✉ yzhang@zuel.edu.cn (Y. Zhang); gyzhou@mail.ccnu.edu.cn (G. Zhou); zwxie@ccnu.edu.cn (Z. Xie); majinjin@mails.ccnu.edu.cn (J. Ma); jhuang@yorku.ca (J.X. Huang)

ORCID(s):

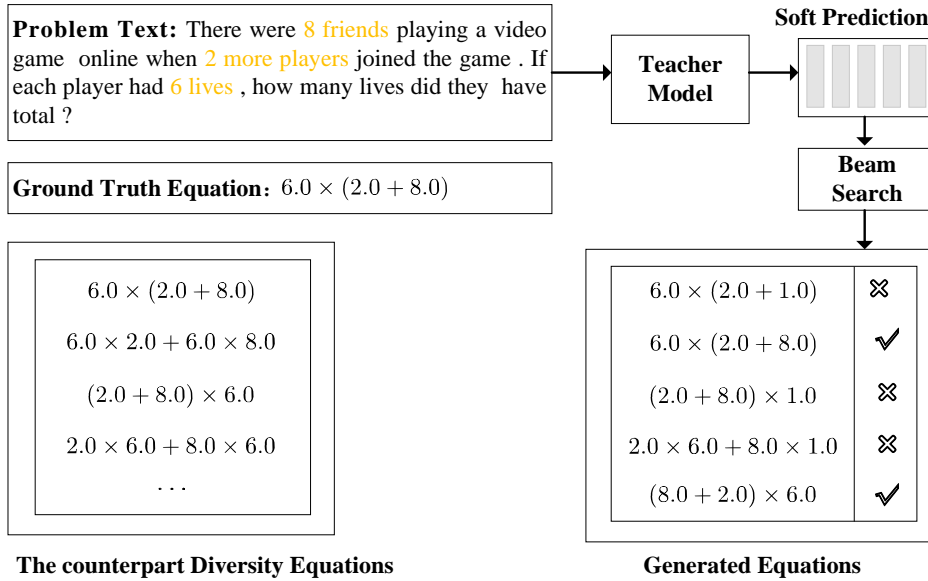


Figure 1: An example of diversity equations and noise soft labels generated by the teacher model Graph2Tree-Z (Zhang et al., 2020b), in which the symbols “x” and “✓” indicate the error and true answer.

equation is “ $6.0 \times (8.0 + 2.0)$ ”, but other equations can also correctly solve this problem, such as “ $6.0 \times 8.0 + 6.0 \times 2.0$ ” and “ $6.0 \times (2.0 + 8.0)$ ”. Conventional Seq2Seq, Seq2Tree, and Graph2Tree models only utilize a single labeled ground-truth equation and a fixed decoder to generate equations, limiting their ability to produce diverse solution equations.

To alleviate these issues, Wang et al. (2018, 2019) employ an equation normalization approach to reduce the diversity of solution equations by imposing restrictions on the order of mathematical operators. However, the effectiveness of this approach has proven to be limited. Zhang et al. (2020a) propose a knowledge distillation approach, called TSN-MD, which encourages the student model to learn diverse equations under the soft supervision of the teacher model by using multiple decoders. However, the use of multiple decoders is less adaptable due to the large and uncertain number of possible equations. Moreover, employing multiple decoders significantly increases inference time, making it impractical for real-world application scenarios. Additionally, some misleading knowledge from the teacher model can impair the student model’s performance. As the example shown in Figure 1, the teacher model generates an incorrect solution equation ranked first in the beam search results. This indicates that the soft prediction of the teacher model follows an inaccurate distribution. Consequently, teaching the student to learn these soft predictions can negatively impact the student model.

Therefore, considering these challenges described above, we propose a novel Diversity-Enhanced Knowledge Distillation (DivKD) model for practical MWP solving. Our approach introduces an Adaptive Diversity Knowledge Distillation (AdaKD) method to effectively transfer the teacher’s diverse knowledge by selectively training the student model with high-quality soft and hard labels. This method adaptively selects the high-quality teacher labels as the diversity target labels for student model learning, and estimates the quality of the intermediate soft labels for training student discernment. The purpose of method is to address the scarcity of diverse solution equations in datasets. Secondly, we design a diversity prior enhanced student model, enabling the model to capture equation diversity by incorporating a Conditional Variational Autoencoder (CVAE) module. In this way, we can model the diversity distribution over possible solution equations by sampling latent variables via a diversity prior network. Extensive experiments conducted on four widely used benchmark MWP datasets demonstrate the effectiveness of our proposed method.

In summary, our main contributions are listed as follows:

- We propose a novel Adaptive Diversity Knowledge Distillation (AdaKD) method, which selectively transfers high-quality knowledge from the pre-trained teacher model to better guide the learning of the student model. This method addresses the scarcity of diverse solution equations in datasets, enabling the student model to acquire both diverse and high-quality information from the pre-trained teacher model.

- We introduce a diversity prior-enhanced student model by incorporating Conditional Variational Autoencoder (CVAE) into existing models, enabling the model to capture the diversity distribution of solution equations. This approach leverages a diversity prior network to sample latent variables, thereby modeling the distribution of possible solution equations and generating diverse solution equations during the testing phase.
- We conducted extensive experiments on four MWP datasets to validate the effectiveness of our proposed Diversity-Enhanced Knowledge Distillation (DivKD) model. The experimental results show that, based on the answer accuracy metric, the proposed methods effectively improve the performance of existing models without sacrificing model efficiency.

The remainder of this paper is organized as follows: Section 2 introduces the related work of this study. Section 3 provides background information and formally defines the research task. Section 4 presents our proposed DivKD approach. Section 5 details the experimental results. Finally, Section 6 concludes the paper and discusses prospects for future work.

2. Related Work

In automatically solving MWPs, deep learning techniques have become the primary approach for MWP solving due to their superior performance and generalization capabilities compared to traditional rule-based and statistical methods (Lu et al., 2023). Recently, the deep learning models (Xie and Sun, 2019; Zhang et al., 2020b, 2022; Bin et al., 2023b; Qin et al., 2023; Bin et al., 2023a) for MWP solving have primarily focused on end-to-end frameworks (e.g., Seq2Seq, Seq2Tree, Graph2Tree). Besides, some subsequent methods, such as PARAMAWPS (Raiyan et al., 2023), DiverseMWP (Zhou et al., 2023) and MathEncoder (Qin et al., 2023), leverage data augmentation, voting, or task-specific pre-training strategies to enhance the performance. However, these approaches are not the focus of this paper due to differing motivations and research goals.

In addition to the aforementioned base solvers, recent pre-trained language models (PLMs) offer opportunity for developing more powerful MWP solvers, including MWP-BERT (Liang et al., 2022), Generate&Rank (Shen et al., 2021) and Deductive Reasoner (Jie et al., 2022). Most recently, the large-scale models have demonstrated remarkable reasoning abilities for MWP solving (Shao et al., 2022; Pi et al., 2022; Liang et al., 2023a). Large language models (LLMs) achieve higher accuracy and explain the solution processes (Touvron et al., 2023; Brown et al., 2020). Large Multimodal models (LMMs) effectively incorporate visual information to aid mathematical reasoning (Shi et al., 2024; Zhuang et al., 2024). However, the computational efficiency makes large size models hard for real-world situations especially in smart education domain where cost and speed are important consideration. Additionally, LLMs such as GPT-3 (Brown et al., 2020) are prone to factual errors (Ouyang et al., 2022), their knowledge bases are not up-to-date. In summary, a computationally efficient solver still needs to be developed to make it practical for real-world situations.

Although the previous MWP solvers have achieved promising performance, they are trained on datasets that provide only a single ground-truth target, which limits the solver's generalization ability. Consequently, some works attempt to construct rich template sketches to accommodate diverse equations during the training process (Mitra and Baral, 2016; Huang et al., 2017; Wang et al., 2019). However, this approach can lead to many different equations with various combinations, resulting in a larger non-deterministic space during the decoding process. Zhang et al. (2020a) propose to extract knowledge from a pre-trained teacher network and encourage multiple student networks to mimic the soft labels generated by the teacher network to learn diverse solutions, named TSN-MD. Subsequent paper (Wang et al., 2022) following TSN-MD focus on unifying the output solution structures through a proposed information storage structure called M-Tree. However, multiple teacher networks or solution structures also require additional model complexity to assess the data quality in mathematical reasoning.

Different from the above-mentioned models, we propose a new framework to better capture the diversity of solution equations in MWPs. Firstly, we design an AdaKD method to select high-quality distillation labels, enriching diversity by incorporating new knowledge from teacher models while reducing the impact of noise information. Secondly, we introduce an enhanced student model with a diversity prior, enabling the model to capture the diversity distribution of solution equations.

Table 1
Notations and Meanings.

Notation	Meanings	Notation	Meanings
x	An MWP textual description.	y	The target mathematical equation.
T	The teacher model in KD.	S	The student model in KD.
$T(x)$	The soft output of teacher model.	$S(x)$	The soft output of student model.
$Encoder(\cdot)$	One of the existing encoder models.	$TreeDecoder(\cdot)$	A Tree-based decoder.
h	The output of $Encoder(x)$.	θ_T	The parameters of the teacher model.
h_y	Representation of target equation.	θ_S	The parameters of the student model.
h_z	The representation of z .	z	Latent distribution over possible solutions.
$\mathcal{N}(\mu, \sigma^2 \mathbf{I})$	Posterior distribution.	$\mathcal{N}(\mu', \sigma'^2 \mathbf{I})$	Diversity prior distribution .
$\mu, \log \sigma^2$	Gaussian parameters for $\mathcal{N}(\mu, \sigma^2 \mathbf{I})$.	$\mu', \log \sigma'^2$	Gaussian parameters for $\mathcal{N}(\mu', \sigma'^2 \mathbf{I})$.
$W_\mu, W_\sigma, b_\mu, b_\sigma$	Trainable parameters for $\{\mu, \log \sigma^2\}$.	$W_{\mu'}, W_{\sigma'}, b_{\mu'}, b_{\sigma'}$	Trainable parameters for $\{\mu', \log \sigma'^2\}$.
\mathcal{L}_{CVAE}	The loss function of CVAE.	$\mathcal{L}_{AdaHardKD}$	The loss function of adaptive hard KD.
w_x	Weight score of teacher's soft labels.	$\mathcal{L}_{AdaSoftKD}$	The loss function of adaptive soft KD.

3. Background

In this section, we introduce the Knowledge Distillation method and the representative teacher-student framework, discuss the relationship in some preliminaries of this work, and finally, summarize the mathematical notations and definitions used in this study in Table 1.

3.1. Knowledge Distillation

Knowledge distillation (KD) (Hinton et al., 2015) has been widely used in deep learning models due to its ability to transfer knowledge from a pre-trained teacher model to a student model. In the process of KD, a teacher model T is pre-trained to generate soft targets, and then a student model S is trained under the supervision of both the ground-truth labels and generated soft labels. Formally, the student model S is trained on a linear combination of two loss functions:

$$\mathcal{L} = (1 - \lambda)\mathcal{L}_{CE} + \lambda\mathcal{L}_{KD}, \quad (1)$$

where λ is a hyper-parameter, \mathcal{L}_{CE} is the cross-entropy between the student output $S(x)$ and the ground-truth label y , which is computed as:

$$\mathcal{L}_{CE} = -y \log \sigma(S(x)). \quad (2)$$

where σ denotes the softmax function. And \mathcal{L}_{KD} is the Kullback-Leibler (KL) divergence loss between student output $S(x)$ and the soft target $T(x)$ generated from teacher model, namely:

$$\mathcal{L}_{KD} = KL(\sigma(S(x)/\tau) || \sigma(T(x)/\tau)) \quad (3)$$

where τ is a temperature hyper-parameter in softmax function.

3.2. Problem Definition

Let $x = \{w_1, w_2, \dots, w_n\}$ denote a math word problem and $y = \{a_1, a_2, \dots, a_m\}$ denote the output solution equation sequence, where w_i in the math word problem is a word token, a_i in the answer equation is a numeric value or a operator, n is the number of words in x and m is the number of elements in y . Given a set of MWPs and corresponding solution equations $D = \{(x, y)\}$, the task of solving math word problems aims to learn a model to map the text sequence of a given math word problem x into an output equation sequence y .

4. Our Method

Figure 2 illustrates the overview of the proposed a diversity-enhanced knowledge distillation (DivKD) model. In our approach, we model the diversity of MWPs using a knowledge distillation framework with a diversity prior. Firstly,

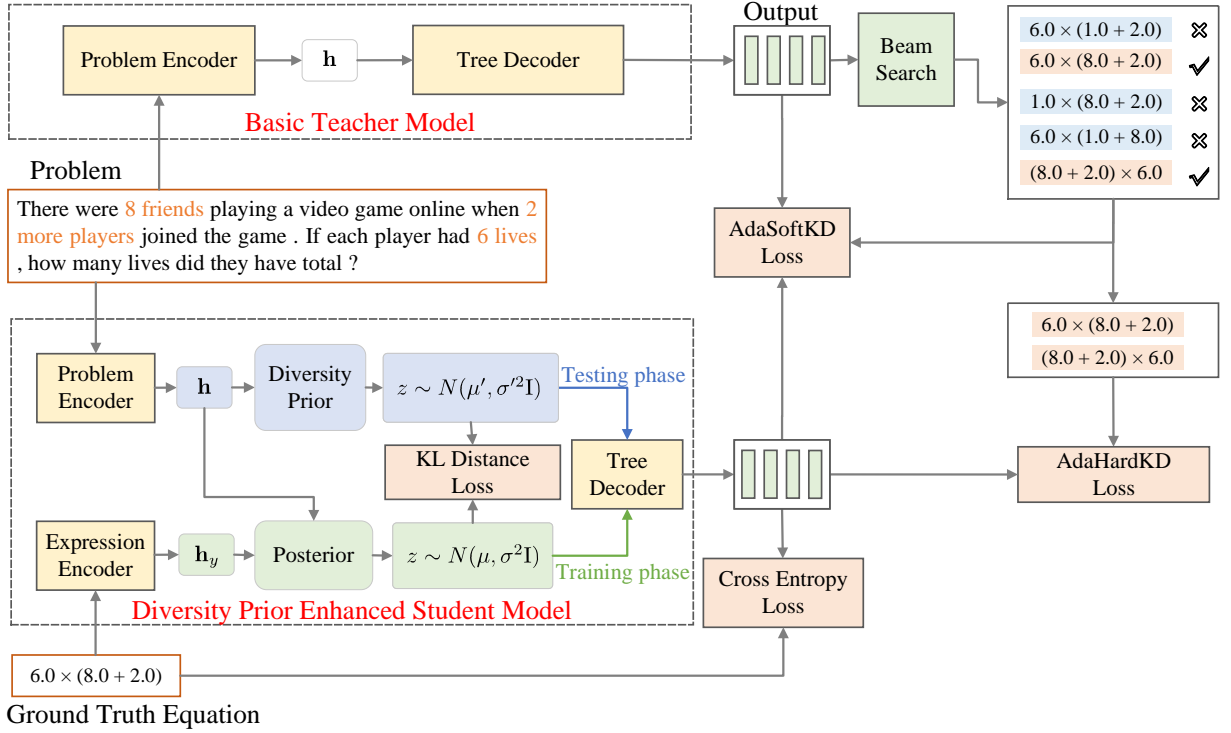


Figure 2: The overview of the proposed DivKD model.

a diversity enhanced teacher-student framework is applied, where a basic teacher is trained to generate new knowledge and a diversity prior enhanced student model is proposed to model the diversity distribution of solution equations. Then, we propose two adaptive knowledge distillation strategies (e.g., adaptive soft distillation and hard distillation) to transfer high-quality knowledge from the teacher model to the student model.

4.1. Diversity Enhanced Knowledge Distillation Framework

In the task of MWP solving, there can be multiple correct solution equations for a given problem, while only one solution equation is annotated in datasets. Most previous models are designed to fit the single annotation solution equation without considering the diversity of solutions. Thanks to the generalization ability of deep learning models, we can leverage KD to learn knowledge beyond the original dataset from a teacher model. The KD framework consists of a teacher model and a student model. In this paper, we select tree-based encoder-decoder models (e.g., GTS (Xie and Sun, 2019) and Graph2Tree-Z (Zhang et al., 2020b)) as basic structures since these models are typical in MWP solving task. GTS (Xie and Sun, 2019) applies a bidirectional gated recurrent unit (BiGRU) as the encoder. Graph2Tree-Z (Zhang et al., 2020b) uses a graph-based encoder (GraphEncoder) to better learn relationships and order information among quantities. Furthermore, we also present a PLM-enhanced version of the Graph2Tree-Z model named Ro-Graph2Tree-Z, which employs a pre-trained language (PLM) model (e.g., RoBERT (Liu et al., 2019)) as an encoder to enhance the understanding of MWPs. It's worth noting that our method is generic and can also be applied to other MWP solvers. In this subsection, we describe the structure of teacher and student models used in our KD framework.

4.1.1. Basic Teacher Model

Firstly, we use the original structure of the encoder-decoder model (e.g., GTS, Graph2Tree-Z and Ro-Graph2Tree-Z) as a teacher model to pre-train on MWP datasets. Given a math problem $x = \{w_1, w_2, \dots, w_n\}$, the teacher model first uses an encoder to model the sequence of input and return the hidden state of the input math problem, which is denoted as:

$$h = \text{Encoder}(x) \quad (4)$$

where h is the output hidden state of the encoder, $Encoder \in \{BiGRU, GraphEncoder, RoBERTaEncoder\}$ is one of the existing text encoder models, and BiGRU, GraphEncoder and RoBERTaEncoder are the encoders in GTS, Graph2Tree-Z and Ro-Graph2Tree-Z, respectively.

Finally, the output feature of the encoder is fed into a tree-based decoder to generate the distribution of solution equations, which is denoted as:

$$p(y|x, \theta_T) = TreeDecoder(h) \quad (5)$$

where θ_T denotes the parameters of the teacher model.

4.1.2. Diversity Prior Enhanced Student Model

The variational autoencoder (VAE) (Kingma and Welling, 2014; Rezende et al., 2014) and its variants, such as conditional variational autoencoder (CVAE) (Sohn et al., 2015), are widely adopted generative models due to their ability of generating diverse data (Zhao et al., 2017; Wu et al., 2020). To model the diversity of solution equations, we introduce a diversity prior to enhancing the student model by incorporating a CVAE component into the original encoder-decoder MWP models (e.g., GTS, Graph2Tree-Z, Ro-Graph2Tree-Z). In the task of MWP solving, the generation of a solution equation y is conditioned on a given math problem x , which can be denoted as $p(y|x)$. Using CVAE, a latent variable z is introduced to capture the latent distribution over possible solution equations and the conditional likelihood calculation is reformulated as:

$$p(y|x) = \int p(y|x, z)p(z|x) \quad (6)$$

Here $p(z|x)$ is a diversity prior distribution which is used to sample diverse latent variable z , and $p(y|x, z)$ is a generative distribution. Thus, the process of generating a solution equation is as follows: (1) sample a latent variable z from the prior distribution $p(z|x)$ and (2) generate solution equation y from the conditional generative distribution $p(y|x, z)$.

It is intractable to directly optimize the conditional log-likelihood of $p(y|x)$ due to the marginalization over the latent variable z . To address this problem, a stochastic gradient variational bayes (SGVB) framework (Kingma and Welling, 2014; Sohn et al., 2015) is applied to train the CVAE model, which converts the original conditional log-likelihood into variational lower bound:

$$\mathcal{L}_{CVAE} = -KL(q(z|x, y)||p(z|x)) + \mathbb{E}_{q(z|x, y)} [\log p(y|x, z)] \leq \log p(y|x) \quad (7)$$

where $q(z|x, y)$ is the posterior distribution. The first term is the KL distance between posterior and prior distribution, and the second term is the cross-entropy between output distribution and ground-truth labels.

Following previous works (Kingma and Welling, 2014; Zhao et al., 2017; Zhang et al., 2016), we assume that z follows a multivariate Gaussian distribution, namely $q(z|x, y) \sim \mathcal{N}(\mu, \sigma^2 \mathbf{I})$ and $p(z|x) \sim \mathcal{N}(\mu', \sigma'^2 \mathbf{I})$. And we use two neural networks to approximate the diversity prior distribution $p(z|x)$ and posterior distribution $q(z|x, y)$. For the posterior distribution $q(z|x, y)$, we first apply a BiGRU network to encode the target equation to obtain the output representation:

$$h_y = BiGRU(y) \quad (8)$$

Then, linear feed-forward networks are applied to obtain the Gaussian parameters μ and $\log \sigma^2$ for the posterior distribution:

$$\begin{aligned} \mu &= W_\mu[h||h_y] + b_\mu \\ \log \sigma^2 &= W_\sigma[h||h_y] + b_\sigma \end{aligned} \quad (9)$$

where W_μ , W_σ , b_μ and b_σ are trainable parameters in neural networks. Similarly, we can compute the Gaussian parameters μ' and $\log \sigma'^2$ for the prior distribution using a linear prior network as follows:

$$\begin{aligned} \mu' &= W_{\mu'}h + b_{\mu'} \\ \log \sigma'^2 &= W_{\sigma'}h + b_{\sigma'} \end{aligned} \quad (10)$$

During training, we sample z from posterior distribution $\mathcal{N}(\mu, \sigma^2 \mathbf{I})$, while the latent variable z is sampled from diversity prior distribution $\mathcal{N}(\mu', \sigma'^2 \mathbf{I})$ during testing. Here, a reparameterization technique (Kingma and Welling, 2014) is applied to obtain the representation of z , namely $h_z = \mu + \sigma \odot \epsilon$, $\epsilon \sim \mathcal{N}(0, \mathbf{I})$. Then, the latent variable can be used to generate solutions via tree-based decoder:

$$p(y|x, \theta_S) = \text{TreeDecoder}(h + h_z) \quad (11)$$

where θ_S is the parameters of the student model.

4.2. Adaptive Diversity Knowledge Distillation

Given that datasets typically provide only a single equation for each math problem, it becomes challenging for a model to gain a comprehensive understanding of diverse solution equations for an MWP. In this study, we employ the KD method to address this issue. Previous KD-based approaches, such as TSN-MD, have solely focused on transferring soft labels from a teacher model to a student model without assessing the quality of the soft labels generated by the teacher. In practice, not all soft labels from the teacher model may enhance the student's performance. Some soft labels may contain noise information that can potentially degrade the student's performance. Therefore, how to measure the quality of teacher knowledge is of great importance to learning a good student model. A simple solution is to measure the quality of soft labels according to the cross-entropy of the teacher's prediction (Li et al., 2021; Wang et al., 2021). However, this method is not well-suited for MWP tasks because there may be multiple correct equations for a single MWP. If we simply employ cross-entropy between soft labels and labelled ground-truth equations, it may result in misjudgment for other unlabeled correct equations. To address this problem, we propose an adaptive diversity knowledge distillation (AdaKD) method for better distillation, which consists of two adaptive KD strategies, namely adaptive hard KD (AdaHardKD) and adaptive soft KD (AdaSoftKD). Specifically, we introduce a beam search-based approach to assess the quality of the teacher's predictions. When provided with the soft predictions from the teacher model, we employ beam search to generate the top- K equations. Subsequently, we calculate the results of these equations and compare them with the ground-truth answer value to determine their correctness. Intuitively, if an equation's result matches the ground-truth value, it is considered a correct equation. Thus, we can obtain a set of correct equations for MWPs, which is denoted as $D^{kd} = \{(x, y^{kd})\}$. These generated equations are viewed as new hard labels to compute the loss of adaptive hard KD, which is calculated as:

$$\mathcal{L}_{AdaHardKD} = \frac{1}{|D^{kd}|} \sum_{(x, y^{kd}) \in D^{kd}} -\log p(y^{kd}|x, \theta_S) \quad (12)$$

In addition to using hard labels, we also introduce an adaptive soft KD strategy to learn high-quality knowledge from the soft labels of the teacher. Intuitively, a better distribution will yield a greater number of correct equations with higher ranks. Consequently, we can assess the quality of the teacher's soft labels based on the outcomes of the beam search. Formally, let $\mathcal{B} = \{(y_k^{kd}, r_k) | 1 \leq k \leq K\}$ denote the results of beam search, where y_k^{kd} is the k -th equation and r_k is its rank in the beam search results. And let $\mathcal{B}' \in \mathcal{B}$ denote the set of correct equations¹. Then, the weight score of the teacher's soft prediction can be computed as:

$$\omega_x = \frac{1}{K} \sum_{(y^{kd}, r) \in \mathcal{B}'} \lambda^r \quad (13)$$

where $\lambda \in (0, 1]$ is an attenuation factor used to assign higher scores to top-ranked equations. By utilizing the weighted scores, we encourage the student model to gain more knowledge from high-quality soft labels. And we define the AdaSoftKD loss function as:

$$\mathcal{L}_{AdaSoftKD} = \frac{1}{|D|} \sum_{(x, y) \in D} \omega_x KL(p(y|x, \theta_S) || q(y|x, \theta_T)) \quad (14)$$

¹To efficiently capture diverse knowledge of expression, we employ answer accuracy rather than expression accuracy as our evaluation metric in this paper. The equation is considered correct in the beam generated by the teacher model when its calculation value equals the ground truth.

Algorithm 1 Diversity-Enhanced Knowledge Distillation

```

1: Input: Training dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ , teacher model  $\mathcal{T}$ , student model  $\mathcal{S}$ , number of beams  $K$ 
2: Output: Trained student model  $\mathcal{S}$ 
3: Initialize teacher model  $\mathcal{T}$  and student model  $\mathcal{S}$ 
4: Pre-train teacher model  $\mathcal{T}$  on  $\mathcal{D}$ 
5: for each batch  $(x, y) \in \mathcal{D}$  do
6:   Obtain teacher predictions  $p(y|x, \theta_T) = \mathcal{T}(x)$ 
7:   Perform beam search to generate top- $K$  equations  $\mathcal{B} = \{(y_k, r_k)\}_{k=1}^K$ 
8:   Identify correct equations  $\mathcal{B}' \subseteq \mathcal{B}$ 
9:   Calculate adaptive hard KD loss:
10:     $D^{kd} \leftarrow \{(x, y^{kd}) \mid y^{kd} \in \mathcal{B}'\}$ 
11:     $\mathcal{L}_{AdaHardKD} \leftarrow -\frac{1}{|D^{kd}|} \sum_{(x, y^{kd}) \in D^{kd}} \log p(y^{kd} | x, \theta_S)$ 
12:   Calculate adaptive soft KD loss:
13:    Compute weight score  $\omega_x$  based on ranks in  $\mathcal{B}$ 
14:     $\mathcal{L}_{AdaSoftKD} \leftarrow \frac{1}{|\mathcal{D}|} \sum_{(x, y) \in \mathcal{D}} \omega_x \cdot KL(p(y|x, \theta_S) || q(y|x, \theta_T))$ 
15:   Total loss:
16:     $\mathcal{L}_{total} = -\mathcal{L}_{CVAE} + \beta \mathcal{L}_{AdaHardKD} + \gamma \mathcal{L}_{AdaSoftKD}$ 
17:   Update student model  $\mathcal{S}$  using optimizer on  $\mathcal{L}_{total}$ 
18: end for
19: Return trained student model  $\mathcal{S}$ 

```

4.3. Training Objective

Finally, the model is trained by using a linear combination of three losses, namely the CVAE loss, the adaptive soft KD loss and the adaptive hard KD loss:

$$\mathcal{L} = -\mathcal{L}_{CVAE} + \beta \mathcal{L}_{AdaHardKD} + \gamma \mathcal{L}_{AdaSoftKD} \quad (15)$$

where β and γ are hyper-parameters controlling the weights of $\mathcal{L}_{AdaHardKD}$ and $\mathcal{L}_{AdaSoftKD}$. To summarize our DivKD method, we provide the algorithm pseudocode in Algorithm 1.

5. Experiments

In this section, we provide a comprehensive description of the datasets employed, the baseline methodologies applied, and the overall experimental settings. Subsequently, we present detailed findings from comparison and ablation experiments, which substantiate the efficacy of our approach. Our focus is directed towards investigating the following five research questions:

- **RQ1:** How does our model perform on the answer accuracy metrics compared to existing MWP solvers?
- **RQ2:** How does our diversity-enhanced knowledge distillation approach, DivKD, perform compared to existing knowledge distillation methods in the MWP task?
- **RQ3:** How does the proposed adaptively diverse knowledge distillation (AdaKD) alleviate the limitations of single labeled ground-truth expression in the dataset?
- **RQ4:** How can our approach improve the performance of existing models without sacrificing model efficiency, especially over other KD methods?
- **RQ5:** How does our approach perform on real-world datasets, particularly in capturing implicit diversity knowledge from problem texts and expressions and alleviating the limitations of single-labeled ground-truth expressions in the datasets?

Table 2
Dataset statistics.

Dataset	#Train	#Dev	#Test	#Avg. Problem Length	#Avg. Operators	Language
Math23K	21,161	1,000	1,000	28.04	2.70	Chinese
MAWPS	1,589	199	199	30.32	1.61	English
MathQA	16,191	2,411	1,605	39.63	4.63	English
SVAMP	3,138	-	1,000	31.87	1.33	English

Table 3
The search scopes and the hyper-parameter setting.

Number of Parameter	Search Scope	Math23K	MAWPS	MathQA	SVAMP
#epoch	{80,100,120,150}	80	80	150	80
batch size	{8,16,30,64}	30	30	16	30
RoBERTa's learning rate	{2e-5,5e-5,5e-6}	5e-5	5e-5	2e-5	5e-5
network layer learning rate	{5e-4,1e-3,5e-3}	1e-3	1e-3	1e-3	1e-3
hidden size	{384,512,768}	512	512	512	512
beam search	{3,5,7}	5	5	5	5
attenuation factor	{0.3,0.5,0.8}	0.8	0.8	0.8	0.8
the weight of AdaHardKD	{0.1,0.2,0.3,0.4}	0.3	0.1	0.1	0.1
the weight of AdaSoftKD	{0.01,0.05,0.1,0.2}	0.1	0.05	0.05	0.05

5.1. General settings

5.1.1. Datasets

In our experiments, we validate our proposed method on four widely used MWP benchmarks: **Math23K** (Wang et al., 2017), **MAWPS** (Koncel-Kedziorski et al., 2016), **MathQA** (Amini et al., 2019) and **SVAMP** (Patel et al., 2021). The overall dataset statistics are presented in Table 2.

Math23K is a large-scale Chinese dataset that consists of 23,161 instances. Each instance is annotated with only one ground-truth equation and its corresponding answer. Following previous studies (Zhang et al., 2020b,a), we adopt the same data splits: 21,161 math problems for training, 1,000 instances for validation, and 1,000 instances for testing. We report results based on answer accuracy on the test data as well as the 5-fold average results.

MAWPS is a standard English MWP-solving dataset that contains 1,987 instances. Since its small size and lack of standard data splits, we perform 5-fold cross-validation and report the average results on this dataset.

MathQA is a large-scale English dataset designed to evaluate the performance of models in solving more complex mathematical problems, where each problem requires multiple operators to derive the final answer.

SVAMP is a challenging English benchmark derived from instances sampled from the MAWPS dataset, consisting of 3,138 training instances and 1,000 test instances. It introduces variations such as noun phrase exchanges and the addition of extra quantities to assess whether NLP models can effectively understand and interpret contextual information.

5.1.2. Evaluation Metric

Answer accuracy and *expression accuracy* are two widely used evaluation metrics in MWP task. *Answer accuracy* indicates that a prediction is correct if the calculated value of the generated expression matches the ground-truth answer. *Expression accuracy* assesses whether the structure of the generated equation is identical to the target expression. The former measures the correctness of the answer, while the latter is concerned with the correctness of the structure. In our paper, we follow most of the prior works (Zhang et al., 2020b; Jie et al., 2022; Bin et al., 2023a) employing the *answer accuracy* as our standard evaluation metrics due to our focus on diversity knowledge inherent in datasets and teacher predictions to enhance the student model's ability to generate diverse equations. *Answer accuracy* highlights the potential of DivKD in diversity-driven mathematical problem-solving.

5.1.3. Experimental Setting

We implement our model using Pytorch² and conduct experiments on Ubuntu 20.04 using a server equipped with two NVIDIA RTX A6000 GPUs, each with 48GB of GPU memory. The Adam optimizer is used in our model, and we use a grid search to select the appropriate hyperparameters based on the evaluation metrics on the test set. For the basic models such as GTS and Graph2Tree-Z, we set the initial learning rate to 1e-3, with the rate halving every 20 epochs. The word embedding dimension is set to 128, and the hidden state dimension to 512. For the basic models with pre-trained language models, such as Ro-Graph2Tree-Z, we use RoBERTa-base (Liu et al., 2019), setting the initial learning rate to 5e-5. For LLaMA (Touvron et al., 2023), we freeze its parameters and use it as an encoder for the input text. To generate more diverse equations, we set the beam search size to 5. The attenuation coefficient λ is set to 0.8, and the weight β for AdaHardKD is set 0.3, 0.1, 0.1 and 0.1 for Math23K, MAWPS, MathQA and SVAMP, respectively. The weight γ for AdaSoftKD is set to 0.1, 0.05, 0.05, 0.05 for the same datasets, respectively. Table 3 lists the main hyperparameters setting in the model and details of their search ranges. Our code and the model's detail are released at <https://github.com/a773938364/DivKD>.

5.1.4. Comparison Models

Based on the architectures employed for math text encoding, we categorize the baseline models into three groups: **Group A**, which includes models utilizing small deep learning architectures such as LSTM, RNN, and GRU; **Group B** comprising models leveraging PLMs like BERT and RoBERTa; and **Group C**, consisting of models incorporating LLMs such as T5 and LLaMA.

- **Basic Models (Group A).** These methods use RNN (e.g., LSTM, GRU) or GNN (e.g., GCN, GAT) to encode the problem text into a feature vector and decompose this vector into an expression. We compare our method with various foundational methods without any language models: DNS (Wang et al., 2017), GTS (Xie and Sun, 2019), T-RNN (Wang et al., 2019), Group-ATT (Li et al., 2019), TSN-MD (Zhang et al., 2020a), SAU-Solver (Qin et al., 2020) and HMS (Lin et al., 2021), Seq2Prog (Amini et al., 2019), NumS2T (Wu et al., 2021), Graph2tree-Z (Zhang et al., 2020b) and MWP-Teacher (Liang and Zhang, 2021).
- **PLM-enhanced Models (Group B).** These methods mainly use pre-trained language models (e.g., BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019)) as encoders and capture the general linguistic and semantic information. These baselines include UniLM (Dong et al., 2019), BERT-Tree (Li et al., 2022), MWP-BERT (Liang et al., 2022), mBERT-LSTM (Tan et al., 2022), SUMC-Solver (Wang et al., 2022), PseDual (Bin et al., 2023b), MWP-NAS (Bin et al., 2023a), C-MWP (Liang et al., 2023b), and PLM-enhanced variants of GTS and Graph2Tree-Z (e.g., Ro-GTS and Ro-Graph2Tree-Z).
- **LLM Models (Group C).** Recently, methods using large models have demonstrated the remarkable reasoning ability for MWP solving. We compare our DivKD with methods using large models listed as follows: T5 (Raffel et al., 2020), GPT-3.5-turbo (Brown et al., 2020), LLaMA (Touvron et al., 2023), Math-PUMA (Zhuang et al., 2024) and Math-LLaVA³ (Shi et al., 2024). We also consider the impact that using more powerful language models (e.g., LLaMA) as teacher models, denoted LLaMA-GTS + DivKD and LLaMA-Graph2Tree-Z + DivKD, which use the LLaMA as the input text encoder.

5.2. Main Results (RQ1 & RQ2)

Table 4 and Table 5 present the experimental results on Math23K, MAWPS, MathQA and SVAMP datasets. We use GTS, Graph2Tree-Z, and their language model variants (e.g., Ro-Graph2Tree-Z, LLaMA-Graph2Tree-Z) as the foundational models within the teacher-student framework, and apply the proposed DivKD method to these basic models. We compare our method against strong baselines and analyze each group.

From Tabel 4 and Table 5, we derive several observations: (1) The models using the graph encoder significantly outperform these models using the sequence encoder (e.g., GTS and Group-ATT vs. Graph2Tree-Z and NumS2T). The reason is that GNN can capture global and long-distance relations in the problem texts, providing more semantic

²<https://pytorch.org/>

³Math-LLaVA is a multimodal LLM designed to exploit visual information to enhance the mathematical reasoning capabilities of multimodal. The comparison of Math-LLaVA is not provided since Math-LLaVA uses picture resources as essential inputs, making it difficult to reproduce their results. We leave these potential improvements in the future work.

Table 4

Answer accuracy on Math23K and MAWPS datasets: Note that the results evaluated on the public test set are denoted as “Math23K”, and the results using 5-fold cross-validation are denoted as “Math23K*” and “MAWPS*”. The symbol “♣” denotes our re-produced results with the public released model.

Category	Model	Size	Math23K	Math23K*	MAWPS*
Group A	DNS	4M	58.1	-	59.5
	T-RNN	8M	66.9	65.2	66.8
	Group-ATT	9M	69.5	66.9	76.1
	GTS	14M	75.6	74.3	82.6
	Graph2Tree-Z	16M	77.4	75.5	83.7
	TSN-MD	22M	77.4	75.1	84.4
	SAU-Solver	21M	76.2	74.8	75.5
	HMS	19M	76.1	-	80.3
	NumS2T	24M	78.1	75.9	84.3
	MWP-Teacher	32M	79.1	77.2	84.2
Group B	UniLM	340M	77.5	-	78.0
	MWP-BERT	130M	84.7	82.4	-
	Ro-GTS	135M	84.0	82.0	87.1
	Ro-Graph2Tree-Z	137M	84.9	82.4	88.7
	SUMC-Solver	135M	77.4	-	79.9
	C-MWP	144M	86.1	84.3	89.1
	PseDual(GCN)	128M	84.6	-	92.4
	MWP-NAS	121M	86.1	-	91.4
Group C	T5	3B	63.2	-	70.4
	GPT-3.5-turbo	175B	54.8	-	91.0
	LLaMA	7B	66.1	-	90.1
	MATH-PUMA	7B	78.5♣	-	82.8♣
	LLaMA-GTS	7B	84.8	82.9	90.6
	LLaMA-Graph2Tree-Z	7B	85.5	84.7	91.3
Ours	GTS + DivKD	14M	77.7 (↑ 2.1)	75.8 (↑ 1.5)	86.2 (↑ 3.6)
	Graph2Tree-Z + DivKD	16M	80.2 (↑ 2.8)	77.6 (↑ 2.1)	86.7 (↑ 4.0)
	Ro-GTS + DivKD	135M	85.1 (↑ 1.1)	83.3 (↑ 1.3)	88.6 (↑ 1.5)
	Ro-Graph2Tree-Z + DivKD	137M	86.2 (↑ 1.3)	84.5 (↑ 2.1)	90.1 (↑ 1.4)
	LLaMA-GTS + DivKD	7B	85.7 (↑ 0.9)	83.6 (↑ 0.7)	91.3 (↑ 0.7)
	LLaMA-Graph2Tree-Z + DivKD	7B	86.7 (↑ 1.2)	85.6 (↑ 0.9)	92.8 (↑ 1.5)

information for decoding. (2) According to the comparison in Group B (e.g., Ro-Graph2Tree-Z vs. MAW-NAS), considering the structure information of expression can significantly improve the performance due to tree structures can decrease the search space and increase the efficiency of decoder. (3) From the results in Group A and Group B, we can find that numerical knowledge can help math operation between numbers, such as UniLM vs. MWP-BERT, achieve better results on answer accuracy metrics. (4) The language model makes it easier to solve MWPs automatically. From Group B and Group A, we observe that models using PLMs achieve significant improvements in solving MWP (e.g., Graph2Tree-Z vs. Ro-Graph2Tree-Z). In particular, we replace the PLMs (e.g., RoBERTa-base) of Graph2Tree-Z with the more powerful language models (e.g., LLaMA), the LLaMA-Graph2Tree-Z achieve the improvement of 2.6% and 1.8% on MAWPS and MathQA. (5) Compared to existing methods, our method learns diverse and high-quality information from the teacher model and sample latent variables, achieving 86.7%, 92.8% and 79.3% on Math23K, 5-fold MAWPS and MathQA. Furthermore, we observe that when using GTS as the basic teacher and student models, the proposed “GTS + DivKD” outperforms the GTS model by 2.1% (77.7% vs. 75.6%) on Math23K, 1.5% (75.8% vs. 74.3%) on 5-fold Math23K, and 3.6% (86.2% vs. 82.6%) on 5-fold MAWPS. When taking Graph2Tree-Z as the base model, our method Graph2Tree-Z + DivKD also achieves better accuracy than the original Graph2Tree-Z across all datasets. These experimental results demonstrate that our proposed DivKD method can effectively improve performance by modelling the diversity of solutions via diversity prior and adaptive diversity KD.

Table 5 presents the accuracy comparison on the more challenging benchmark datasets, MathQA and SVAMP. We can find that these basic models using our DivKD method achieve significant improvements. For instance, our solver

Table 5

Results of various models on MathQA and SVAMP benchmarks: †Results are from (Bin et al., 2023b). The symbol “♣” denotes our re-produced results with the public released model.

Category	Model	Size	MathQA	SVAMP
Group A	DNS	4M	65.7	18.7
	Group-ATT	9M	70.4	21.5
	Seq2Prog	10M	57.2	-
	NumS2T	24M	72.7	37.0
Group B	BERT-Tree	121M	75.1	32.4
	mBERT-LSTM	132M	74.7	-
	PseDual(GCN)	128M	78.9	-
	Ro-GTS	135M	73.5	41.0
	Ro-Graph2Tree-Z	137M	74.4	43.8
Group C	GPT-3.5-turbo	175B	73.5	69.3 [†]
	LLaMA	7B	75.7	38.0
	Math-PUMA	7B	54.3 [♣]	68.2 [♣]
	LLaMA-Graph2Tree-Z	7B	76.2	51.3
Ours	Ro-GTS + DivKD	135M	77.1 († 3.6)	42.8 († 1.8)
	Ro-Graph2Tree-Z + DivKD	137M	78.7 († 4.3)	45.3 († 1.5)
	LLaMA-Graph2Tree-Z + DivKD	7B	79.3 († 3.1)	52.5 († 1.2)

“Ro-Graph2Tree-Z + DivKD” surpasses Ro-Graph2Tree-Z by 4.3% (78.7% vs. 74.4%) on MathQA, highlighting the effectiveness of our approach in enhancing problem-solving capabilities. We also analyze the improvements, and the results further confirm that our DivKD approach generate diverse and correct solutions, while the baselines are limited by their single solution approach. On the SVAMP dataset, our approach achieves an answer accuracy of 52.5%, which is lower than GPT-3.5-turbo and Math-PUMA, but higher than other baseline methods. We observe that the average number of operators in the SVAMP is 1.33, significantly lower than in MathQA (average of 4.63) and Math23K (average of 2.70), indicating that SVAMP solutions are less complex. This complexity disparity explains why large language models (LLMs) perform well on SVAMP. In contrast, our DivKD approach is designed to model solution diversity more effectively, which is crucial for addressing more complex problem sets.

Among the baselines, TSN-MD is a related model that also uses the KD method with GTS as basic teacher and student models. From Table 4, we can see that our method “GTS + DivKD” outperforms TSN-MD on Math23K, 5-fold Math23K and 5-fold MAWPS (e.g., 77.7% vs. 77.4% on Math23k, 75.8% vs. 75.1% on 5-fold Math23K and 86.2% vs. 84.4% on 5-fold MAWPS). TSN-MD simply performs distillation using soft labels from the teacher model without considering the quality of these soft labels, which could result in learning misinformation due to noisy labels. Instead, our method is able to adaptively select high-quality labels for KD, thereby mitigating the impact of noisy soft labels generated from the teacher model. Additionally, TSN-MD uses multiple decoders in the student model to generate diverse solution equations, which is time-consuming. In contrast, we propose a diversity prior to efficiently modeling the diversity of solution equations using only one decoder, without increasing the complexity of the original GTS model. Therefore, our proposed DivKD method is superior to the previous KD-based method (e.g., TSN-MD) in both performance and efficiency.

5.3. Ablation Study (RQ3)

To better understand the impact of different components in the proposed DivKD method, we conduct ablation studies by constructing some variants of DivKD. As shown in Table 6, we use GTS and Graph2Tree-Z as basic teacher and student models and design four variants with different KD strategies:

- “GTS/Graph2Tree-Z + DivKD w/o S&H” denotes the GTS/Graph2Tree-Z model enhanced with diversity prior using CVAE, but without using adaptive soft and hard KD method.
- “GTS/Graph2Tree-Z + DivKD w/o H” only uses AdaSoftKD to train the student model without using AdaHardKD.

Table 6
Ablation Study.

Model	Math23K	MAWPS
GTS	75.6	82.6
GTS + DivKD	77.7	86.2
GTS + DivKD w/o S&H	76.4	84.6
GTS + DivKD w/o H	76.9	85.0
GTS + DivKD w/o S	77.2	85.9
GTS + CVAE + SoftKD	76.6	84.7
Graph2Tree-Z	77.4	83.7
Graph2Tree-Z + DivKD	80.2	86.7
Graph2Tree-Z + DivKD w/o S&H	78.2	85.1
Graph2Tree-Z + DivKD w/o H	78.9	85.5
Graph2Tree-Z + DivKD w/o S	79.5	86.1
Graph2Tree-Z + CVAE + SoftKD	78.4	85.2

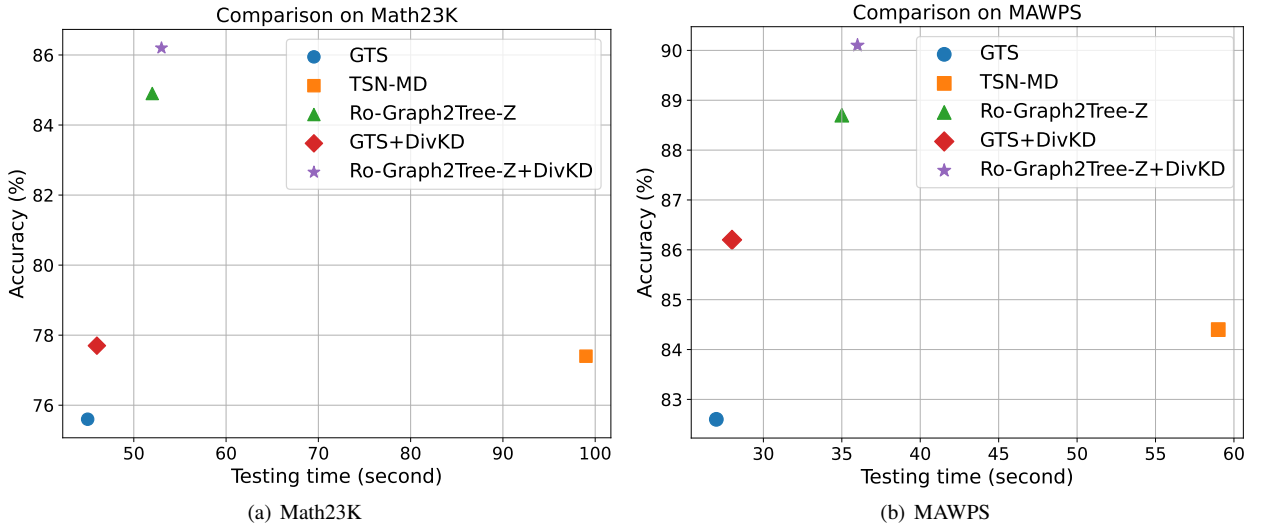


Figure 3: Comparison of testing time between baselines and our proposed methods (e.g., GTS+DivKD and Ro-Graph2TreeZ+DivKD) for Math23K and MAWPS datasets.

- “GTS/Graph2Tree-Z + DivKD w/o S” only uses AdaHardKD to train the student model without using AdaSoftKD.
- “GTS/Graph2Tree-Z + CVAE + SoftKD” denotes using a conventional soft KD method to train the student model, without considering the quality of soft labels.

As shown in Table 6, the proposed DivKD model without using KD outperforms the basic GTS/Graph2Tree-Z models on both datasets. For example, on MAWPS dataset, the “GTS + DivKD w/o S&H” achieves an accuracy of 84.6%, surpassing the performance of GTS (e.g., 82.6%), and “Graph2Tree-Z + DivKD w/o S&H” also outperforms original Graph2Tree-Z by a margin of 1.4% (e.g., 85.1% vs. 83.7%). This reveals that our DivKD model can effectively learn the diversity of solution equations by using diversity prior enhanced model.

In addition, by applying KD, the performance can be further improved. The “GTS/Graph2Tree-Z +DivKD w/o H” models, which adaptively distill knowledge from high-quality soft labels, obtain better performance than “GTS/Graph2Tree-Z + CVAE + SoftKD”, indicating the effectiveness of the proposed adaptive KD method. Furthermore, we can observe that “GTS/Graph2Tree-Z + DivKD w/o S” performs better than other KD methods (e.g., conventional soft KD and adaptive soft KD), showing that the adaptive hard KD is superior to other KD methods since we can directly evaluate the answer of hard equations to filter out noise labels. Using both the adaptive hard KD and

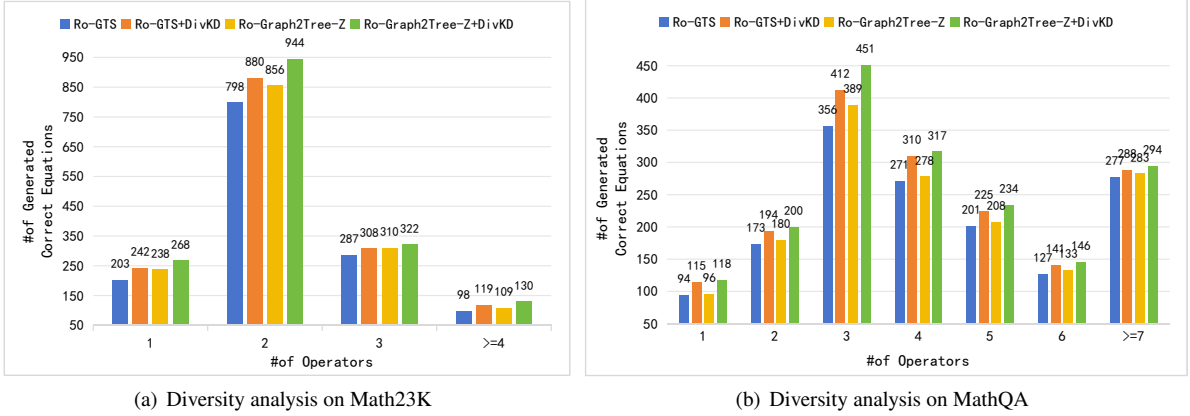


Figure 4: A quantitative analysis of generated correct expressions between basic models (e.g., Ro-GTS) and the student models (e.g., Ro-GTS+DivKD) on Math23K and MathQA datasets.

adaptive soft KD, our “GTS + DivKD” can achieve the best performance. In summary, the ablation studies demonstrate that all the components in the proposed DivKD model play important roles in the improvement of performance.

5.4. Efficiency Analysis (RQ4)

In this section, we compare the efficiency of the proposed model with existing baseline models, including GTS, TSN-MD and Ro-Graph2Tree-Z. GTS serves as the foundational model structure for both TSN-MD and our proposed “GTS+DivKD”. TSN-MD is a related model that aims to model the diversity of solution equations using multiple decoders. Ro-Graph2Tree-Z and “Ro-Graph2Tree-Z+DivKD” utilize pre-trained Roberta model and graph structure to enhance the encoding of input problem. We evaluate these models on the test sets of Math23K and MAWPS and present their testing times in Figure 3. All experiments are conducted on a single RTX A6000 48G graphics card to ensure a fair comparison.

Comparing among GTS, TSN-MD and “GTS+DivKD”, we observe that the running time of “GTS+DivKD” is similar to that of the basic GTS model, whereas TSN-MD requires significantly more time for inference than both GTS and “GTS+DivKD”. This increased inference time in TSN-MD is due to its use of two decoders to generate solution equations. In contrast, “GTS+DivKD” utilizes only one decoder and incorporates a CVAE to better model the diversity of solution equations, thereby enhancing performance without compromising efficiency. Furthermore, by employing PLMs, Ro-Graph2Tree-Z and “Ro-Graph2Tree-Z+DivKD” improve performance on both Math23K and MAWPS datasets with only a slight increase in computational time compared to GTS and “GTS+DivKD”. But the testing time of Ro-Graph2Tree-Z and “Ro-Graph2Tree-Z+DivKD” are still much less than TSN-MD model. This is because the tree-based decoding process is very time-consuming, and the TSN-MD model, which uses multiple decoders, will inevitably significantly increase the testing time. Therefore, using multiple decoders to improve the performance is not a good choice. Different from TSN-MD, our DivKD method can be conveniently plugged into different models and can improve model effectiveness without causing an increase in inference time. This renders our approach feasible for practical application.

5.5. Quantitative Analysis (RQ5)

In this section, we conduct a quantitative analysis to count the number of the correct solution equations generated by the teacher models (e.g., Ro-GTS and Ro-Graph2Tree-Z) and student models (e.g., “Ro-GTS + DivKD” and “Ro-Graph2Tree-Z + DivKD”) in Figure 4. On the Math23K and MathQA datasets, we can see that the student models produce a significantly higher number of correct equations compared to the teacher models. For instance, on the subset of MathQA with 3 operators, our “Ro-Graph2Tree-Z+DivKD” outperforms the Ro-Graph2Tree-Z by 62. This improvement is attributed to DivKD’s ability to learn the diversity distribution of solution equations, thereby generating a wider range of potential solutions during the test phase. Instead, Ro-GTS and Ro-Graph2Tree-Z are limited to learning only one ground-truth solution equation provided for each math problem in the benchmark datasets. Overall, these

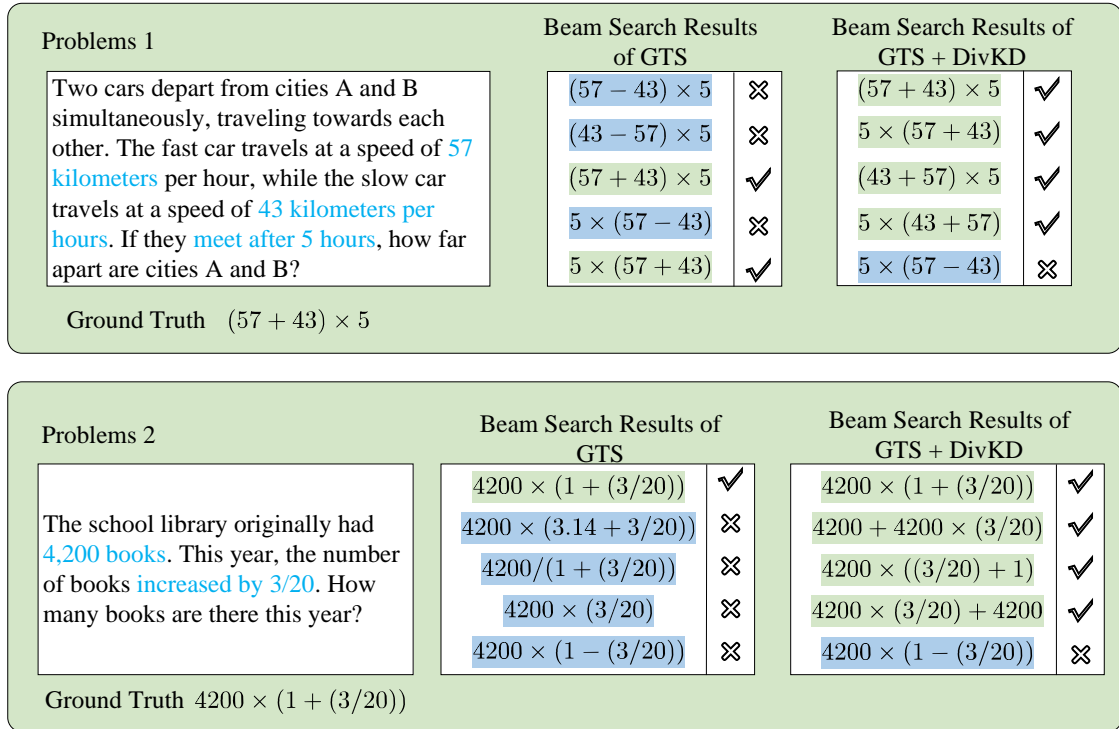


Figure 5: Two examples of the generated results on Math23K using GTS and the proposed GTS+DivKD.

results demonstrate that our method enhances both the diversity and correctness of generated equations, showcasing DivKD's capability to produce more accurate and varied solutions.

5.6. Case Study (RQ5)

To better understand the superiority of the proposed method, we conduct a case study by comparing the results generated by the original GTS and our proposed "GTS+DivKD". As shown in Figure 5, for Problem 1, the original GTS model generates a wrong solution equation at the first rank of the five beam search results and only two correct solution equations are ranked at the top five results. Instead, our proposed "GTS+DivKD" can correctly predict four solution equations and rank them at top positions. In the second example, although the GTS predicts a correct solution equation ranking at the first position, it fails to model the diversity of solution equations since only one correct equation is ranked in the top 5 results. Different from the original GTS, our proposed "GTS+DivKD" can generate diverse correct solution equations in the beam search results. These examples demonstrate that the proposed DivKD method is able to model the diversity distribution of multiple solution equations for MWPs.

5.7. Discussion

In recent years, LLMs and LMMs such as GPT-3.5-turbo (Brown et al., 2020) and Math-PUMA (Zhuang et al., 2024) have made remarkable strides in solving MWPs with the help of chain-of-thought (CoT) prompting (Chen et al., 2022; Wang et al., 2023; Xie et al., 2024). Although our model does not outperform these advanced large-scale models on the MWP task, we believe our work contributes significant value to the research field. Real-time or resource-constrained is crucial for many real-world applications, such as smart education. However, existing LLMs have hundreds of millions of parameters and exorbitant API, resulting in computational inefficiency and impracticality for real-world situations. In contrast, our method achieves competitive performance compared to LLMs while offering faster response times and requiring fewer parameters for the MWP task. Furthermore, our method is a natural solution that distills the knowledge from a bigger model into smaller, more efficient student models, which can potentially mitigate the issues of LLMs for practical situations. Objectively speaking, our goal is designing practical MWP solvers to balance effectiveness and efficiency.

6. Conclusion and Future Work

In this paper, introduce a novel Diversity-Enhanced Knowledge Distillation (DivKD) model for solving Math Word Problems (MWP). Our approach effectively models diverse solution equations in MWPs and extracts high-quality knowledge from the teacher model. Specifically, we enhance the student model with a diversity prior by integrating a Conditional Variational Autoencoder (CVAE) module into conventional encoder-decoder architectures (e.g., GTS, Graph2Tree-Z, and Ro-Graph2Tree-Z). This integration enables the model to capture the diversity of equations. Moreover, to overcome the limited annotation equations in datasets, we propose an Adaptive Diversity Knowledge Distillation (AdaKD) method that selects high-quality knowledge from the teacher model, providing an additional supervisory signal to guide the student model's learning process. Empirical evaluations on four widely used benchmark datasets demonstrate the superiority of the proposed DivKD model in MWP.

Future research can proceed in two primary directions. First, despite the strong performance of large-scale models such as GPT-3.5-turbo and LLaMA in MWP solving, their practical deployment remains limited due to their substantial model sizes and high computational costs. Online educational applications require models that can serve millions of students with low latency while designing personalized learning paths, presenting significant challenges in both effectiveness and efficiency. Therefore, a promising direction is to focus on compressing these large models into smaller, more efficient versions without compromising performance, thereby facilitating their integration into personalized mobile educational platforms. Second, we aim to evaluate the robustness of the proposed DivKD method on more complex MWP datasets, including CM17K (Qin et al., 2021), GSM8K (Cobbe et al., 2021), and MATH (Hendrycks et al., 2021). Assessing performance on these datasets will provide deeper insights into the model's ability to handle a wider variety of problem complexities and further validate its effectiveness.

Data Availability

Data will be made available on request.

ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China under Grant 62377021, the China Postdoctoral Science Foundation under Grant Number 2024M751062, financially supported by self-determined research funds of CCNU from the colleges' basic research and operation of MOE (No. CCNU24XJ010, CCNU22QN015 and CCNU24ai011), the Natural Science Foundation of Hubei Province for Distinguished Young Scholars (No. 2023AFA096), and the Wuhan Knowledge Innovation Project (No. 2022010801010278). This work was also supported by the Natural Sciences and Engineering Research Council (NSERC) of Canada, an NSERC CREATE award in ADERSIM3 and the York Research Chairs (YRC) program.

References

- Amini, A., Gabriel, S., Lin, S., Koncel-Kedziorski, R., Choi, Y., Hajishirzi, H., 2019. MathQA: Towards interpretable math word problem solving with operation-based formalisms, in: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Association for Computational Linguistics, Minneapolis, Minnesota. pp. 2357–2367. URL: <https://aclanthology.org/N19-1245>, doi:10.18653/v1/N19-1245.
- Bin, Y., Han, M., Shi, W., Wang, L., Yang, Y., Ng, S.K., Shen, H., 2023a. Non-autoregressive math word problem solver with unified tree structure, in: Bouamor, H., Pino, J., Bali, K. (Eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Singapore. pp. 3290–3301. URL: <https://aclanthology.org/2023.emnlp-main.199>, doi:10.18653/v1/2023.emnlp-main.199.
- Bin, Y., SHI, W., Ding, Y., Yang, Y., Ng, S.K., 2023b. Solving math word problems with reexamination. URL: <https://openreview.net/forum?id=4c6s9is9DV>.
- Bobrow, D., 1964. Natural language input for a computer problem solving system. Technical report, Massachusetts Institute of Technology .
- Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D., 2020. Language models are few-shot learners, in: Proceedings of the 34th International Conference on Neural Information Processing Systems, Curran Associates Inc., Red Hook, NY, USA. pp. 1877 – 1901.
- Chen, W., Ma, X., Wang, X., Cohen, W.W., 2022. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. URL: <https://doi.org/10.48550/arXiv.2211.12588>, arXiv:2211.12588.

- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., Schulman, J., 2021. Training verifiers to solve math word problems. URL: <https://arxiv.org/abs/2110.14168>, arXiv:2110.14168.
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2019. BERT: Pre-training of deep bidirectional transformers for language understanding, in: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Association for Computational Linguistics, Minneapolis, Minnesota. pp. 4171–4186. URL: <https://aclanthology.org/N19-1423>, doi:10.18653/v1/N19-1423.
- Dong, L., Yang, N., Wang, W., Wei, F., Liu, X., Wang, Y., Gao, J., Zhou, M., Hon, H., 2019. Unified language model pre-training for natural language understanding and generation, in: Wallach, H.M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E.B., Garnett, R. (Eds.), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, pp. 13042–13054. URL: <https://dl.acm.org/doi/10.5555/3454287.3455457>.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., Steinhardt, J., 2021. Measuring mathematical problem solving with the math dataset. URL: https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/be83ab3ecd0db773eb2dc1b0a17836a1-Paper-round2.pdf.
- Hinton, G.E., Vinyals, O., Dean, J., 2015. Distilling the knowledge in a neural network. URL: <http://arxiv.org/abs/1503.02531>, arXiv:1503.02531.
- Hosseini, M.J., Hajishirzi, H., Etzioni, O., Kushman, N., 2014. Learning to solve arithmetic word problems with verb categorization, in: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, Doha, Qatar. pp. 523–533. doi:10.3115/v1/D14-1058.
- Huang, D., Shi, S., Lin, C., Yin, J., 2017. Learning fine-grained expressions to solve math word problems, in: Palmer, M., Hwa, R., Riedel, S. (Eds.), Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017, Association for Computational Linguistics. pp. 805–814. doi:10.18653/v1/d17-1084.
- Jie, Z., Li, J., Lu, W., 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), Association for Computational Linguistics, Dublin, Ireland. pp. 5944–5955. URL: <https://aclanthology.org/2022.acl-long.410>, doi:10.18653/v1/2022.acl-long.410.
- Kingma, D.P., Welling, M., 2014. Auto-encoding variational bayes. URL: <http://arxiv.org/abs/1312.6114>.
- Koncel-Kedziorski, R., Hajishirzi, H., Sabharwal, A., Etzioni, O., Ang, S.D., 2015. Parsing algebraic word problems into equations. Transactions of the Association for Computational Linguistics 3, 585–597. URL: <https://aclanthology.org/Q15-1042>, doi:10.1162/tacl_a_00160.
- Koncel-Kedziorski, R., Roy, S., Amini, A., Kushman, N., Hajishirzi, H., 2016. Mawps: A math word problem repository, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, San Diego, California. pp. 1152–1157. URL: <https://aclanthology.org/N16-1136>, doi:10.18653/v1/N16-1136.
- Kushman, N., Zettlemoyer, L., Barzilay, R., Artzi, Y., 2014. Learning to automatically solve algebra word problems, in: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, Volume 1: Long Papers, The Association for Computer Linguistics. pp. 271–281. doi:10.3115/v1/p14-1026.
- Li, J., Wang, L., Zhang, J., Wang, Y., Dai, B.T., Zhang, D., 2019. Modeling intra-relation in math word problems with different functional multi-head attentions, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Florence, Italy. pp. 6162–6167. URL: <https://aclanthology.org/P19-1619>, doi:10.18653/v1/P19-1619.
- Li, L., Lin, Y., Ren, S., Li, P., Zhou, J., Sun, X., 2021. Dynamic knowledge distillation for pre-trained language models, in: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Online and Punta Cana, Dominican Republic. pp. 379–389. URL: <https://aclanthology.org/2021.emnlp-main.31>, doi:10.18653/v1/2021.emnlp-main.31.
- Li, Z., Zhang, W., Yan, C., Zhou, Q., Li, C., Liu, H., Cao, Y., 2022. Seeking patterns, not just memorizing procedures: Contrastive learning for solving math word problems, in: Findings of the Association for Computational Linguistics: ACL 2022, Association for Computational Linguistics, Dublin, Ireland. pp. 2486–2496. URL: <https://aclanthology.org/2022.findings-acl.195>, doi:10.18653/v1/2022.findings-acl.195.
- Liang, Z., Yu, W., Rajpurohit, T., Clark, P., Zhang, X., Kalyan, A., 2023a. Let GPT be a math tutor: Teaching math word problem solvers with customized exercise generation, in: Bouamor, H., Pino, J., Bali, K. (Eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Singapore. pp. 14384–14396. URL: <https://aclanthology.org/2023.emnlp-main.889/>, doi:10.18653/v1/2023.emnlp-main.889.
- Liang, Z., Zhang, J., Guo, K., Wu, X., Shao, J., Zhang, X., 2023b. Compositional mathematical encoding for math word problems, in: Rogers, A., Boyd-Graber, J., Okazaki, N. (Eds.), Findings of the Association for Computational Linguistics: ACL 2023, Association for Computational Linguistics, Toronto, Canada. pp. 10008–10017. URL: <https://aclanthology.org/2023.findings-acl.635>, doi:10.18653/v1/2023.findings-acl.635.
- Liang, Z., Zhang, J., Wang, L., Qin, W., Lan, Y., Shao, J., Zhang, X., 2022. MWP-BERT: Numeracy-augmented pre-training for math word problem solving, in: Findings of the Association for Computational Linguistics: NAACL 2022, Association for Computational Linguistics, Seattle, United States. pp. 997–1009. URL: <https://aclanthology.org/2022.findings-naacl.74>, doi:10.18653/v1/2022.findings-naacl.74.
- Liang, Z., Zhang, X., 2021. Solving math word problems with teacher supervision, in: Zhou, Z.H. (Ed.), Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21, International Joint Conferences on Artificial Intelligence Organization. pp. 3522–3528. URL: <https://doi.org/10.24963/ijcai.2021/485>, doi:10.24963/ijcai.2021/485. main Track.
- Lin, X., Huang, Z., Zhao, H., Chen, E., Liu, Q., Wang, H., Wang, S., 2021. HMS: A hierarchical solver with dependency-enhanced understanding for math word problem, in: Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI Press. pp. 4232–4240.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V., 2019. Roberta: a robustly optimized bert pretraining approach. URL: <http://arxiv.org/abs/1907.11692>, arXiv:1907.11692.

- Lu, P., Qiu, L., Yu, W., Welleck, S., Chang, K.W., 2023. A survey of deep learning for mathematical reasoning, in: Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Toronto, Canada. pp. 14605–14631. URL: <https://aclanthology.org/2023.acl-long.817>.
- Mitra, A., Baral, C., 2016. Learning to use formulas to solve simple arithmetic problems, in: Erk, K., Smith, N.A. (Eds.), Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Berlin, Germany. pp. 2144–2153. URL: <https://aclanthology.org/P16-1202/>, doi:10.18653/v1/P16-1202.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P.F., Leike, J., Lowe, R., 2022. Training language models to follow instructions with human feedback, in: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (Eds.), Advances in Neural Information Processing Systems, Curran Associates, Inc.. pp. 27730–27744. URL: https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf.
- Patel, A., Bhattamishra, S., Goyal, N., 2021. Are NLP models really able to solve simple math word problems?, in: Toutanova, K., Rumshisky, A., Zettlemoyer, L., Hakkani-Tur, D., Beltagy, I., Bethard, S., Cotterell, R., Chakraborty, T., Zhou, Y. (Eds.), Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, Online. pp. 2080–2094. URL: <https://aclanthology.org/2021.naacl-main.168>, doi:10.18653/v1/2021.naacl-main.168.
- Pi, X., Liu, Q., Chen, B., Ziyadi, M., Lin, Z., Fu, Q., Gao, Y., Lou, J.G., Chen, W., 2022. Reasoning like program executors, in: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Abu Dhabi, United Arab Emirates. pp. 761–779. URL: <https://aclanthology.org/2022.emnlp-main.48>.
- Qin, J., Liang, X., Hong, Y., Tang, J., Lin, L., 2021. Neural-symbolic solver for math word problems with auxiliary tasks, in: Zong, C., Xia, F., Li, W., Navigli, R. (Eds.), Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, Online. pp. 5870–5881. URL: <https://aclanthology.org/2021.acl-long.456>, doi:10.18653/v1/2021.acl-long.456.
- Qin, J., Lin, L., Liang, X., Zhang, R., Lin, L., 2020. Semantically-aligned universal tree-structured solver for math word problems, in: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, Online. pp. 3780–3789. URL: <https://aclanthology.org/2020.emnlp-main.309>, doi:10.18653/v1/2020.emnlp-main.309.
- Qin, J., Yang, Z., Chen, J., Liang, X., Lin, L., 2023. Template-based contrastive distillation pretraining for math word problem solving. IEEE Transactions on Neural Networks and Learning Systems, 1–13doi:10.1109/TNNLS.2023.3265173.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J., 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. J. Mach. Learn. Res. 21. URL: <https://dl.acm.org/doi/10.5555/3455716.3455856>.
- Raiyan, S.R., Faiyaz, M.N., Kabir, S.M.J., Kabir, M., Mahmud, H., Hasan, M.K., 2023. Math word problem solving by generating linguistic variants of problem statements, in: Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 4: Student Research Workshop), Association for Computational Linguistics, Toronto, Canada. pp. 362–378. URL: <https://aclanthology.org/2023.acl-srw.49>, doi:10.18653/v1/2023.acl-srw.49.
- Rezende, D.J., Mohamed, S., Wierstra, D., 2014. Stochastic backpropagation and approximate inference in deep generative models, in: Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21–26 June 2014, JMLR.org. pp. 1278–1286. URL: <http://proceedings.mlr.press/v32/rezende14.html>.
- Roy, S., Roth, D., 2018. Mapping to declarative knowledge for word problem solving. Transactions of the Association for Computational Linguistics 6, 159–172. URL: <https://aclanthology.org/Q18-1012>, doi:10.1162/tacl_a_00012.
- Shao, Z., Huang, F., Huang, M., 2022. Chaining simultaneous thoughts for numerical reasoning, in: Findings of the Association for Computational Linguistics: EMNLP 2022, Association for Computational Linguistics, Abu Dhabi, United Arab Emirates. pp. 2533–2547. URL: <https://aclanthology.org/2022.findings-emnlp.187>.
- Shen, J., Yin, Y., Li, L., Shang, L., Jiang, X., Zhang, M., Liu, Q., 2021. Generate & rank: A multi-task framework for math word problems, in: Findings of the Association for Computational Linguistics: EMNLP 2021, Association for Computational Linguistics, Punta Cana, Dominican Republic. pp. 2269–2279. URL: <https://aclanthology.org/2021.findings-emnlp.195>, doi:10.18653/v1/2021.findings-emnlp.195.
- Shi, S., Wang, Y., Lin, C.Y., Liu, X., Rui, Y., 2015. Automatically solving number word problems by semantic parsing and reasoning, in: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. pp. 1132–1142.
- Shi, W., Hu, Z., Bin, Y., Liu, J., Yang, Y., Ng, S.K., Bing, L., Lee, R.K.W., 2024. Math-llava: Bootstrapping mathematical reasoning for multimodal large language models. URL: <https://arxiv.org/abs/2406.17294>, arXiv:2406.17294.
- Sohn, K., Lee, H., Yan, X., 2015. Learning structured output representation using deep conditional generative models, in: Cortes, C., Lawrence, N.D., Lee, D.D., Sugiyama, M., Garnett, R. (Eds.), Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, pp. 3483–3491. URL: <https://proceedings.neurips.cc/paper/2015/hash/8d55a249e6baa5c06772297520da2051-Abstract.html>.
- Tan, M., Wang, L., Jiang, L., Jiang, J., 2022. Investigating math word problems using pretrained multilingual language models, in: Proceedings of the 1st Workshop on Mathematical Natural Language Processing (MathNLP), Association for Computational Linguistics, Abu Dhabi, United Arab Emirates (Hybrid). pp. 7–16. URL: <https://aclanthology.org/2022.mathnlp-1.2>.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., Lample, G., 2023. Llama: Open and efficient foundation language models. URL: <https://arxiv.org/abs/2302.13971>, arXiv:2302.13971.
- Wang, B., Ju, J., Fan, Y., Dai, X., Huang, S., Chen, J., 2022. Structure-unified M-tree coding solver for math word problem, in: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Abu Dhabi, United Arab Emirates. pp. 8122–8132. URL: <https://aclanthology.org/2022.emnlp-main.556>.

- Wang, F., Yan, J., Meng, F., Zhou, J., 2021. Selective knowledge distillation for neural machine translation, in: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, Online. pp. 6456–6466. URL: <https://aclanthology.org/2021.acl-long.504>, doi:10.18653/v1/2021.acl-long.504.
- Wang, L., Wang, Y., Cai, D., Zhang, D., Liu, X., 2018. Translating a math word problem to a expression tree, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Brussels, Belgium. pp. 1064–1069. URL: <https://aclanthology.org/D18-1132>, doi:10.18653/v1/D18-1132.
- Wang, L., Zhang, D., Zhang, J., Xu, X., Gao, L., Dai, B.T., Shen, H.T., 2019. Template-based math word problem solvers with recursive neural networks, in: Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, AAAI Press, Honolulu, Hawaii, USA. pp. 7144 – 7151. URL: <https://doi.org/10.1609/aaai.v33i01.33017144>, doi:10.1609/aaai.v33i01.33017144.
- Wang, X., Wei, J., Schuurmans, D., Le, Q.V., Chi, E.H., Narang, S., Chowdhery, A., Zhou, D., 2023. Self-consistency improves chain of thought reasoning in language models. URL: <https://openreview.net/pdf?id=1PL1NIMMrw>.
- Wang, Y., Liu, X., Shi, S., 2017. Deep neural solver for math word problems, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark. pp. 845–854. URL: <https://aclanthology.org/D17-1088>, doi:10.18653/v1/D17-1088.
- Wu, C., Wang, P.Z., Wang, W.Y., 2020. On the encoder-decoder incompatibility in variational text modeling and beyond, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Online. pp. 3449–3464. URL: <https://aclanthology.org/2020.acl-main.316>, doi:10.18653/v1/2020.acl-main.316.
- Wu, Q., Zhang, Q., Wei, Z., Huang, X., 2021. Math word problem solving with explicit numerical values, in: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, Association for Computational Linguistics, Online. pp. 5859–5869. URL: <https://aclanthology.org/2021.acl-long.455>, doi:10.18653/v1/2021.acl-long.455.
- Xiao, J., Huang, L., Song, Y., Tang, N., 2023. A recursive tree-structured neural network with goal forgetting and information aggregation for solving math word problems. Information Processing & Management 60, 103324. URL: <https://www.sciencedirect.com/science/article/pii/S0306457323000614>, doi:<https://doi.org/10.1016/j.ipm.2023.103324>.
- Xie, R., Huang, C., Wang, J., Dhingra, B., 2024. Adversarial math word problem generation. URL: <https://arxiv.org/abs/2402.17916>, arXiv:2402.17916.
- Xie, Z., Sun, S., 2019. A goal-driven tree-structured neural model for math word problems, in: Kraus, S. (Ed.), Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019, ijcai.org. pp. 5299–5305. URL: <https://doi.org/10.24963/ijcai.2019/736>, doi:10.24963/ijcai.2019/736.
- Zhang, B., Xiong, D., Su, J., Duan, H., Zhang, M., 2016. Variational neural machine translation, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Austin, Texas. pp. 521–530. URL: <https://aclanthology.org/D16-1050>, doi:10.18653/v1/D16-1050.
- Zhang, J., Lee, R.K.W., Lim, E.P., Qin, W., Wang, L., Shao, J., Sun, Q., 2020a. Teacher-student networks with multiple decoders for solving math word problem, in: Bessiere, C. (Ed.), Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20, International Joint Conferences on Artificial Intelligence Organization. pp. 4011–4017. URL: <https://doi.org/10.24963/ijcai.2020/555>, doi:10.24963/ijcai.2020/555. main track.
- Zhang, J., Wang, L., Lee, R.K.W., Bin, Y., Wang, Y., Shao, J., Lim, E.P., 2020b. Graph-to-tree learning for solving math word problems, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Online. pp. 3928–3937. URL: <https://aclanthology.org/2020.acl-main.362>, doi:10.18653/v1/2020.acl-main.362.
- Zhang, W., Shen, Y., Ma, Y., Cheng, X., Tan, Z., Nong, Q., Lu, W., 2022. Multi-view reasoning: Consistent contrastive learning for math word problem, in: Findings of the Association for Computational Linguistics: EMNLP 2022, Association for Computational Linguistics, Abu Dhabi, United Arab Emirates. pp. 1103–1116. URL: <https://aclanthology.org/2022.findings-emnlp.79>.
- Zhang, Y., Zhou, G., Xie, Z., Huang, J.X., 2024. Number-enhanced representation with hierarchical recursive tree decoding for math word problem solving. Information Processing & Management 61, 103585. URL: <https://www.sciencedirect.com/science/article/pii/S0306457323003229>, doi:<https://doi.org/10.1016/j.ipm.2023.103585>.
- Zhao, T., Zhao, R., Eskénazi, M., 2017. Learning discourse-level diversity for neural dialog models using conditional variational autoencoders, in: Barzilay, R., Kan, M. (Eds.), Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Association for Computational Linguistics. pp. 654–664. URL: <https://doi.org/10.18653/v1/P17-1061>, doi:10.18653/v1/P17-1061.
- Zhou, Z., Ning, M., Wang, Q., Yao, J., Wang, W., Huang, X., Huang, K., 2023. Learning by analogy: Diverse questions generation in math word problem, in: Findings of the Association for Computational Linguistics: ACL 2023, Association for Computational Linguistics, Toronto, Canada. pp. 11091–11104. URL: <https://aclanthology.org/2023.findings-acl.705>, doi:10.18653/v1/2023.findings-acl.705.
- Zhuang, W., Huang, X., Zhang, X., Zeng, J., 2024. Math-puma: Progressive upward multimodal alignment to enhance mathematical reasoning. URL: <https://arxiv.org/abs/2408.08640>, arXiv:2408.08640.