

# Expression Syntax Information Bottleneck for Math Word Problems

Jing Xiong  
xiongj69@mail2.sysu.edu.cn  
Sun Yat-sen University  
Shenzhen, Guangdong, China

Chengming Li\*  
lichengming@mail.sysu.edu.cn  
Sun Yat-sen University  
Shenzhen, Guangdong, China

Min Yang\*  
min.yang@siat.ac.cn  
SIAT, Chinese Academy of Sciences  
Shenzhen, Guangdong, China

Xiping Hu  
huxiping@mail.sysu.edu.cn  
Sun Yat-sen University  
Shenzhen, Guangdong, China

Bin Hu  
bh@lzu.edu.cn  
Lanzhou University  
Lanzhou, Gansu, China

## ABSTRACT

Math Word Problems (MWP) aims to automatically solve mathematical questions given in texts. Previous studies tend to design complex models to capture additional information in the original text so as to enable the model to gain more comprehensive features. In this paper, we turn our attention in the opposite direction, and work on how to discard redundant features containing spurious correlations for MWP. To this end, we design an **Expression Syntax Information Bottleneck** method for MWP (called *ESIB*) based on variational information bottleneck, which extracts essential features of expression syntax tree while filtering latent-specific redundancy containing syntax-irrelevant features. The key idea of *ESIB* is to encourage multiple models to predict the same expression syntax tree for different problem representations of the same problem by mutual learning so as to capture consistent information of expression syntax tree and discard latent-specific redundancy. To improve the generalization ability of the model and generate more diverse expressions, we design a self-distillation loss to encourage the model to rely more on the expression syntax information in the latent space. Experimental results on two large-scale benchmarks show that our model not only achieves state-of-the-art results but also generates more diverse solutions. The code is available.<sup>1</sup>

## KEYWORDS

Math Word Problems, Mutual learning, Spurious correlations, Variational information bottleneck

### ACM Reference Format:

Jing Xiong, Chengming Li, Min Yang, Xiping Hu, and Bin Hu. 2022. Expression Syntax Information Bottleneck for Math Word Problems. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*, July 11–15, 2022, Madrid, Spain. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/3477495.3531824>

\*Corresponding authors.

<sup>1</sup>[https://github.com/menik1126/math\\_ESIB](https://github.com/menik1126/math_ESIB)

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

*SIGIR '22, July 11–15, 2022, Madrid, Spain*

© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-8732-3/22/07...\$15.00

<https://doi.org/10.1145/3477495.3531824>

## 1 INTRODUCTION

Math Word Problems (MWP) is challenging and draws much attention from researchers in the field of natural language processing [3, 25] and information retrieval (e.g., mathematical understanding) [8, 26]. MWP aims to automatically answer mathematical questions given in a natural language, which requires the model not only understand what facts are presented in a text, but also possess the reasoning capability to answer the mathematical question. Table 1 shows three examples of MWP with three mathematical problems and their solution expressions with answer.

Inspired by the success of deep learning [14, 21], the encoder-decoder framework with attention mechanisms [2] have been dominated in MWP [18–20], which bring the state-of-the-art to a new level. The key idea is to use an encoder to learn representations of problem text and employ a decoder to generate the corresponding solution expression and answer. Subsequently, several studies propose sequence-to-tree models, which explore the tree structure information presented in the text and improve the generation of solution expressions [23, 24, 28].

However, the previous MWP methods appear to rely on spurious correlations between the shallow heuristics in problem and solution expression [15]. For example, as shown in Table 1, previous models may associate Problem 1 and Problem 2 with the mathematical formula “ $x \times y \div z$ ”, since these two problems have similar semantic patterns, e.g., *calculating the speed*. Based on this association, the models could generate wrong solution expression for Problem 3 which has similar semantic problem expression like the text segment “place A to place B” in Problems 1-2. In particular, the models that learn spurious correlations are more likely to generate wrong solution expression “ $220 \times 25\% \div 30\%$ ”, rather than “ $220 \div (25\% + 30\%)$ ” for Problem 3. We define such a false association as spurious correlation. Some recent studies have revealed that MWP solvers relying on spurious correlations could achieve high accuracy [9, 15]. These models can even compute correct answers without paying attention to the question part in the problem such as the text segment “*how many kilometers is the total length of the two places?*” in Problem 3 *calculating the distance*. In addition, the solution expression is sensitive to the perturbed latent representations [10], since the semantically similar mathematical problems, even with totally different solution expressions and questions, can be encoded closely in the latent space. We believe it is an evidence that redundant information containing spurious correlations is encoded

**Table 1: Three examples of MWP.**

<b>Problem 1:</b> From place A to place B, if a bicycle travels 16 kilometers per hour, it can be reached in 4 hours. If it only takes 2 hours by car, how many kilometers per hour does the car travel?
<b>Solution Expression 1:</b> $16 \times 4 \div 2$ <b>Answer:</b> 32
<b>Problem 2:</b> Uncle Jack drove from place A to place B, and it took 6 hours to arrive at the speed of 70 kilometers per hour. When I returned, I accelerated the speed due to the task. It only took 4 hours to return to the first place. What was the speed when I returned?
<b>Solution Expression 2:</b> $70 \times 6 \div 4$ <b>Answer:</b> 105
<b>Problem 3:</b> A car travels 25% of the whole journey from place A to place B in the first hour, 30% of the whole journey in the second hour, a total of 220 kilometers in two hours, how many kilometers is the total length of the two places?
<b>Solution Expression 3:</b> $220 \div (25\% + 30\%)$ <b>Answer:</b> 400
<b>Wrong Solution Expression 3:</b> $220 \times 25\% \div 30\%$

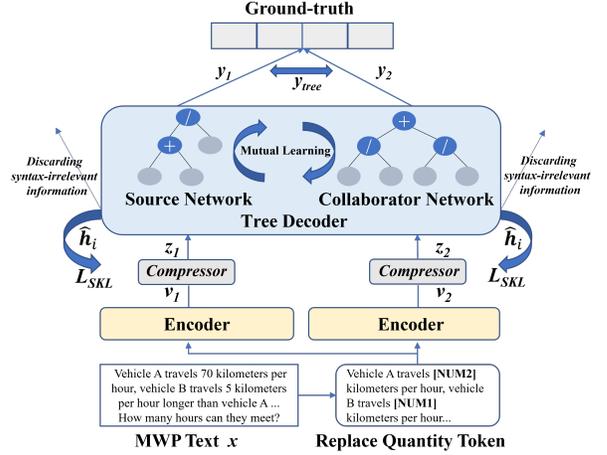
in the latent representation. Therefore, it is necessary to alleviate the spurious correlations by compressing the latent representations for math expressions while filtering latent-specific redundancy.

To solve the above challenges, we design a **Expression Syntax Information Bottleneck** method for MWP based on variational information bottleneck (VIB) [1], which aims to discard redundant information containing spurious correlations [7, 13]. Our key idea is to encourage multiple models to predict the same math expression with different problem representations of the same problem so as to capture consistent information about expression syntax tree in expressions and discard latent-specific redundancy containing syntax-irrelevant information. In addition, we leverage mutual learning [29] for learning variational information bottleneck, which can effectively reduce the latent-specific redundancy. Inspired by the observation that there are usually multiple solutions to solve a problem, we also design a self-distillation loss which encourages the decoder to rely on the syntax information in latent space, enabling the model to generate diverse solutions.

We summarize our main contributions as follows. (1) We are the first to reduce spurious correlations so as to improve the performance of MWP. (2) We propose a novel expression syntax information bottleneck method for MWP, which extracts essential syntax information of math expression and filters redundant information containing spurious correlations. (3) We design a self-distillation loss to encourage the model to generate more diverse solution expressions. (4) Extensive experiments on two benchmark datasets show that our model outperforms the strong baselines in a noticeable margin.

## 2 METHODOLOGY

A math word problems (MWP) can be denoted by a projection  $F : x \mapsto z \mapsto y$ , where  $x = \{w_1, w_2, \dots, w_m\}$  is the problem sequence with  $m$  words,  $z$  is the compressed representation of the original problem  $x$  and  $y = \{o_1, o_2, \dots, o_n\}$  is the solution expression of the problem with  $n$  words. The goal of MWP is to establish a model  $F$  which generates a correct solution expression  $y$  and calculates the correct answer for the problem  $x$ .



**Figure 1: Overview of the proposed method ESIB.**

As illustrated in Figure 1, the proposed *ESIB* is composed of a source network (denoted as SN) and a collaborator network (denoted as CN). The two networks are optimised collaboratively and capture consistent information across different problem representations throughout the training process. We use the deep variational information bottleneck (VIB) framework as the backbone of both SN and CN. The VIB aims to generate problem representation  $z \in \mathbb{R}^s$  by compressing and discarding redundant information in latent representation  $v \in \mathbb{R}^d$  without reducing essential information related to the target  $y$ , where  $d$  denotes hidden size of decoder and  $s$  denotes the dimension of problem representation. So the above projection  $F$  can be rewritten as  $F : x \mapsto v \mapsto z \mapsto y$ , where  $v$  denotes latent representation and  $z$  denotes problem representation sampled from  $e^\mu(v) \sim N(e^\mu(v), e^\sigma(v))$  ( $e$  denotes a dense layer).

### 2.1 Encoder-Compressor-Decoder Architecture

**Encoder.** We adopt the RoBERTa model [11] as our encoder. We pass the problem sequence  $x$  into the RoBERTa model and obtain latent representation  $v$  ([CLS] vector output by the pre-trained language model (PLM) and its dimension is converted from 768 to  $d$  by a dense layer). In order to model the relationship between the quantities in the PLM, we set up a learnable quantity embedding matrix  $Q_t = \{q_1, q_2, \dots, q_n\}$ , similar to the learnable position embedding in BERT [5]. When passing the sequence  $x$  into the encoder, we replace the each quantity in the sequence  $x$  (i.e., the numbers in the problem) with a embedding  $q_i \in Q_t$ .

**Compressor.** Our model takes the encoder output  $v$  and feeds it into a variational information bottleneck [1] module which outputs a sampled vector  $z$ . This part removes redundant information by optimizing a variational upper bound  $\mathcal{V}_{IB}$ . We will detail how to optimize  $\mathcal{V}_{IB}$  in Section 2.2.

**Decoder.** Our decoder follows the GTS model [24]. We use the sampled representation  $z$  to initialize the initial state of the decoder and the recursive operation of the decoder to construct  $y$  by the order of pre-order traversal. First, the root node  $q_{root}$  (middle operator part) is first generated. Then, we generate the left child

node  $q_l$ . Finally, we generate the right child node  $q_r$ . This process has been iterated until the leaf nodes are generated. The attention mechanism is applied to learn the global context vector  $G_i$  which is utilized to generate the current node token  $\hat{y}_i$ . Here we denote the digital embedding after being encoded by the encoder as  $Q$ . Mathematically, we define the attention mechanism as follows:

$$G_i = \begin{cases} \text{Attention}(x, q_{\text{root}}, q_l), & q_l \notin \emptyset. \\ \text{Attention}(x, q_{\text{root}}, q_{sl}), & q_{sl} \notin \emptyset. \\ \text{Attention}(x, q_{\text{root}}), & q_l, q_{sl} \in \emptyset. \end{cases} \quad (1)$$

$$\hat{y}_i, \hat{h}_i = \text{Predict}(G_i, Q), \quad (2)$$

where  $\text{Predict}(\cdot)$  is prediction layer for producing tree nodes  $\hat{y}_i$  and hidden state of decoding step  $i$  (denoted as  $\hat{h}_i \in \mathbb{R}^d$ ).

If the current node is an operator, we will generate the left and right child nodes and push them into the stack in the tree decoder according to the top-down method. If it is a number, we will perform the merge operation until the leaf nodes in the stack pop out, and the result of the merge is pushed into the left child node stack for attention operation. The merge operation will pop the required node  $q_{\text{op}}$  and  $q_{\text{subtree}}$  from an embedding stack. Finally, the merge operation outputs the answer of the mathematical expression. This recursive construction process can be defined as follows:

$$q_l = \text{Left}(\hat{G}, \hat{y}_i, q_{\text{root}}). \quad (3)$$

$$q_r = \text{Right}(\hat{G}, \hat{y}_i, q_{\text{root}}). \quad (4)$$

$$q_m = \text{Merge}(q_{\text{op}}, q_{\text{subtree}}, q_{m-1}). \quad (5)$$

## 2.2 Information Bottleneck

Our *ESIB* method is driven by Information Bottleneck (IB) [1, 17] that forces features to be concise by filtering the task-irrelevant information (i.e., syntax-irrelevant spurious correlations for MWP). Specifically, suppose we know the joint distribution  $p(x, y)$ , our goal is to learn a problem representation  $z$  that maximizes its predictive power for generating  $y$ , subject to the constraints of the amount of information it carries about  $x$ :

$$\mathcal{L}_{IB} = I(y; z) - \lambda I(x; z), \quad (6)$$

where  $I(\cdot; \cdot)$  denotes the mutual information.  $\lambda$  is a Lagrangian multiplier that controls the trade-off between the sufficiency (the task performance quantified by  $I(y; z)$ ) and the minimality (the complexity of the representation quantified by  $I(x; z)$ ). In this paper, we focus on compressing the redundancy in latent representation  $v$  (denoted as a substitute for the problem  $x$ ). Following [1], Equation (6) can be variationally upper bounded by:

$$\mathcal{V}_{IB} = \frac{1}{N} \sum_{n=1}^N \{ \lambda D_{\text{KL}}(e(z | v^n) \| b(z)) - \mathbb{E}_e \log d(y^n | e(z^n, \epsilon)) \}, \quad (7)$$

where  $N$  is the number of data,  $e(z^n, \epsilon)$  transforms the  $z$  into initial state of the decoder (denoted as  $\bar{z} \in \mathbb{R}^{m \cdot d}$ ) and  $e(z | v^n)$  transforms the representation  $v$  into two tensors:  $e^\mu(v)$  denotes the features-mean and  $e^\sigma(v)$  denotes the diagonal features-covariance. We use the reparameterization to obtain the compressed representation  $z = e(v, \epsilon) = e^\mu(v) + \epsilon e^\sigma(v)$  with  $\epsilon \sim N(0, 1)$ . The prior distribution of  $z$  (denoted as  $b(z)$ ) is a standard normal distribution. The decoder  $d$  converts the sampled representation  $z$  into a mathematical expression  $y$ , and calculates the answer.

With the compression capability of VIB, it is possible to lose the necessary information about the target  $y$  in the feature  $z$  when optimizing the trade-off between the compression and redundancy. We will introduce how to identify the syntax-irrelevant information in Section 2.3.

## 2.3 Latent-specific Redundancy

In this section, we demonstrate that the latent-specific redundancy containing syntax-irrelevant information can be effectively reduced by using mutual learning [29]. Inspired by the method forcing two networks to learn from each other [29], we encourage two networks to produce compressed representation  $z$  filtered latent-specific redundancy for keeping all predictive information about expression syntax trees in expression  $y$  by passing an representation of expression syntax tree (denoted as  $y_{\text{tree}}$  obtained by concatenating each  $\hat{h}_i$  defined in equation(2) when decoder predicts the expression syntax tree of the math expression) to each other.

We attempt to learn a vector  $z$  that contains expression syntax tree information about solution expression  $y$  as much as possible and achieve this goal by optimizing the mutual information of  $v$  and  $z$ . We take CN as an example, and factorize the mutual information [6] between  $v_1$  and  $z_1$  as follows:

$$I(v_1; z_1) = \underbrace{I(v_1; z_1 | v_2)}_{\text{Latent-specific Redundancy}} + \underbrace{I(v_2; z_1)}_{\text{Consistent Information}}, \quad (8)$$

where  $v_1$  and  $z_1$  denote latent and problem representation for CN respectively;  $v_2$  and  $z_2$  denote that of SN.  $I(v_1; z_1 | v_2)$  indicates that the information contained in  $z_1$  is unique to  $v_1$  and cannot be inferred by representation  $v_2$  [6]. We call  $I(v_1; z_1 | v_2)$  as latent-specific redundancy or syntax-irrelevant information. It can be discarded by minimizing  $I(v_1; z_1 | v_2)$  [6] which can be upper bounded by the following inequality:

$$\begin{aligned} I(v_1; z_1 | v_2) &= \mathbb{E}_{v_1, v_2 \sim S_1(v|x)} \mathbb{E}_{z_1, z_2 \sim S_2(z|v)} \left[ \log \frac{p(z_1 | v_1)}{p(z_1 | v_2)} \right] \\ &= \mathbb{E}_{v_1, v_2 \sim S_1(v|x)} \mathbb{E}_{z_1, z_2 \sim S_2(z|v)} \left[ \log \frac{p(z_1 | v_1) p(z_2 | v_2)}{p(z_2 | v_2) p(z_1 | v_2)} \right] \\ &= D_{\text{KL}}[p(z_1 | v_1) \| p(z_2 | v_2)] - D_{\text{KL}}[p(z_2 | v_1) \| p(z_2 | v_2)] \\ &\leq D_{\text{KL}}[p(z_1 | v_1) \| p(z_2 | v_2)]. \end{aligned} \quad (9)$$

Inspired by [6], we approximate the upper bound above by replacing  $z$  with  $y_{\text{tree}}$ . Since  $y_{\text{tree}}$  generated from  $z$  contains all the information of the representation  $z$  and all the latent-specific redundancy to be discarded. In addition, the parameters to be optimized about  $z$  are included in the decoder. Considering the above points, we utilize  $D_{\text{KL}}(\mathbb{P}_{z_1} \| \mathbb{P}_{z_2})$  ( $\mathbb{P}_{z_1}$  denotes  $p(y_{\text{tree}_1} | z_1)$  and  $\mathbb{P}_{z_2}$  denotes  $p(y_{\text{tree}_2} | z_2)$  right) as an upper bound to approximate  $D_{\text{KL}}[p(z_1 | v_1) \| p(z_2 | v_2)]$ . Similarly, we can use  $D_{\text{KL}}(\mathbb{P}_{z_2} \| \mathbb{P}_{z_1})$  to minimize  $I(v_2; z_2 | v_1)$  for SN. We introduce the objective  $\mathcal{L}_{SKL}$  to minimize the latent-specific redundancy for both  $z_1$  and  $z_2$ :

$$\mathcal{L}_{SKL} = \min_{\theta, \phi} \mathbb{E}_{v_1, v_2 \sim E_\theta(v|x)} \mathbb{E}_{z_1, z_2 \sim E_\phi(z|v)} [D_{SKL}[\mathbb{P}_{z_1} \| \mathbb{P}_{z_2}]] \quad (10)$$

where  $\theta$  and  $\phi$  denote the parameters of CN and SN. The two networks are optimized alternately during training.  $\mathbb{P}_{z_1} = p_\theta(y_1 | z_1)$  and  $\mathbb{P}_{z_2} = p_\phi(y_2 | z_2)$  denote the concatenation of the output distributions at each step of the model prediction of CN and SN

respectively.  $D_{SKL}$  denotes symmetrized KL divergence obtained by averaging the expected value of  $D_{KL}(\mathbb{P}_{z_1} \parallel \mathbb{P}_{z_2})$  and  $D_{KL}(\mathbb{P}_{z_2} \parallel \mathbb{P}_{z_1})$ . We calculate this loss by mutual learning [29] between CN and SN. In the mutual learning setup, in an iteration, the model will compute  $\mathcal{L}_{SKL_1}$  and  $\mathcal{L}_{SKL_2}$  for CN and SN respectively. In addition, we need to maximize  $I(v_2; z_1)$  to ensure that the compressed representation  $z_1$  has enough information to predict  $y$ . We use the chain rule to decompose  $I(v_2; z_1)$  into the following two terms:

$$I(v_2; z_1) = \underbrace{I(v_2; z_1 | y)}_{\text{Redundancy}} + \underbrace{I(z_1; y)}_{\text{Predictive Information}}, \quad (11)$$

where  $y$  represents ground-truth solution expression. In practice, we can maximize  $I(z_1; y) = \mathbb{E}_\epsilon \log d(y^n | e(x^n, \epsilon))$  (included in  $\mathcal{V}_{IB}$ ) which is calculated to compress redundant information  $I(v_2; z_1 | y)$  and indirectly maximize  $I(v_2; z_1)$ . Ideally,  $I(v_2; z_1 | y)$  should be zero. As suggested in [6], we minimize latent-specific redundancy by jointly minimizing  $I(v_1; z_1 | v_2)$  and maximizing  $I(v_2; z_1)$ . We define the redundancy terms in Eq. (8) and Eq. (11) as the syntax-irrelevant information.

## 2.4 Self-distillation Loss

In this section, we introduce a novel self-distillation loss to increase the diversity of generated expressions. Suggested by [7], the  $I(v; z | y)$  in conditional information bottleneck can be variationally upper bounded by:

$$D_{KL}(e(z | v^n) \| b(z | y^n)), \quad (12)$$

where  $e(z | v^n)$  defined in equation (2) is moved towards the conditional marginal  $b^\mu(y) \sim N(b^\mu(y), b^\sigma(y))$ . We modify equation (12) as:

$$\mathcal{V}_{SDL} = \frac{1}{N} \sum_{n=1}^N D_{SKL}(\bar{y} \| \bar{z}), \quad (13)$$

where  $\bar{y}$  denotes that averaging the all  $\hat{h}_i$ . Intuitively,  $\mathcal{V}_{SDL}$  make the decoder more rely on latent space of  $z$  which contains expression syntax tree information for all expressions when predicting expression  $y$ . Benefiting from the randomness of  $z$ , the model can generate more diverse solution expressions. Finally, we calculate the loss functions  $\mathcal{L}_1$  for SN and  $\mathcal{L}_2$  for CN as follows:

$$\mathcal{L}_1 = \mathcal{V}_{IB_1} + \mathcal{V}_{SDL_1} + \alpha \times \mathcal{L}_{SKL_1}. \quad (14)$$

$$\mathcal{L}_2 = \mathcal{V}_{IB_2} + \mathcal{V}_{SDL_2} + \alpha \times \mathcal{L}_{SKL_2}. \quad (15)$$

where  $\alpha$  is a proportional coefficient.  $\mathcal{V}_{IB_1}, \mathcal{V}_{SDL_1}, \mathcal{L}_{SKL_1}$  are the training objectives for SN.  $\mathcal{V}_{IB_2}, \mathcal{V}_{SDL_2}, \mathcal{L}_{SKL_2}$  are the training objectives for CN. Based on empirical observation, although  $\mathcal{V}_{SDL}$  can increase the diversity of solution expressions, it also reduces accuracy of the final answer.

## 3 EXPERIMENTAL SETUP

*Datasets.* We conduct experiments on two benchmark MWP datasets: Math23k [22] and Ape210k [30]. Math23k contains 22162/1000 questions for training/testing, respectively. Ape210k is composed of 166,270 questions for training, 4,157 questions for validation, and 4,159 questions for testing.

*Implementation Details.* The word embedding size of decoder is set to 1024 and proportional coefficient  $\alpha$  in loss function is set to 0.005. We set the dimension of vectors  $z$  to 50. When the encoder is BERT, we set the dimension of  $z$  to 32. We adopt RoBERTa [11] as the problem encoder. Following RoBERTa’s setting, the hidden size of the encoder is set to 768, and we set the hidden size of the decoder to 1024. We used Adamw [12] as the optimizer with the learning rate as 5e-5. The mini-batch size is set to 16. We adopt a beam search with the size of 5. Dropout (dropout rate = 0.5) is employed to avoid overfitting. For Ape210K, we set the maximum sequence length of questions as 150 and that of solution expressions as 50, similar to [23]. Our model takes 80 epochs on Math23k and 50 epochs on Ape210k for convergence.

*Baselines.* We compare our model with several strong baseline methods, including NumS2T [23], TSN-MD [27], MATH-EN [18], Multi-E/D [16], GTS [24], StackDecoder [4], KAS2T [23], Graph2tree [28], and Ape [30].

## 4 EXPERIMENTAL RESULTS

### 4.1 Main Results

The evaluation metric is answer accuracy. Table 2 show the performance comparison of our model with baseline methods on Math23K and Ape210k, respectively. In Table 2, since there is a trade-off between the variety of expressions and the correctness of the answer, we do not add  $\mathcal{V}_{SDL}$  into the source network (SN) and the collaborator network (CN). From Table 2 we can observe that our models (both CN and SN) achieve consistently and substantially better performance than the compared methods. In addition, the accuracy of CN is higher than that of SN. This is because, in one iteration, CN is provided with  $y_{tree}$  predicted by SN when SN has not been trained by ground-truth then SN is provided with  $y_{tree}$  predicted by CN when CN has been trained by ground-truth.

We also measure the accuracy of solution expression. We consider a solution expression as correct when the predicted expression exactly matches the ground truth solution. Generally, the mathematical expression generated by the tree decoder conforms to the syntactic specification. As long as the answer obtained by the expression operation is consistent with the ground-truth, then we consider the expression to be a valid solution. We take the subtraction value between answer accuracy (denoted as Answer-Acc) and solution expression accuracy (denoted as Expression-Acc) as the diversity evaluation metric (denoted as Diversity) of the generated solution expressions. As shown in Table 3, our model can generate more diverse solution expressions that are not included in ground-truth expressions. As expected, the model with  $\mathcal{V}_{SDL}$  has better diversity but lower answer accuracy.

### 4.2 Ablation Study

We conduct ablation test on Math23k to analyze the impact of different components in *ESIB*. Since the best results are produced by CN, we only conduct ablation study on CN. First, we remove the mutual learning, denoted as CN w/o MT. Second, we remove the VIB from SN and CN to evaluate the impact of VIB (denoted as CN w/o VIB). In addition, we report the results by removing both MT and VIB (denoted as CN w/o MT+VIB). To evaluate the impact of

**Table 2: Solution accuracy of *ESIB* and various baselines. Note that Math23K denotes results on public test set.**

Models	Math23k	Ape210k
StackDecoder	-	52.2
GTS	75.6	67.7
KAS2T	76.3	68.7
TSN-MD	77.4	-
Graph2Tree	77.4	-
Ape	-	70.2
NumS2T	78.1	70.5
Multi-E/D	78.4	-
SN	<b>84.2</b>	<b>76.3</b>
CN	<b>85.9</b>	<b>76.8</b>

**Table 3: Accuracy of diversity generation on Math23k with Self-distillation Loss.**

Models	Answer-Acc	Equation-Acc	Diversity
MATH-EN	66.7	60.1	6.6
GTS	75.6	64.8	10.8
TSN-MD	77.4	65.8	11.6
CN with $\mathcal{V}_{SDL}$	<b>85.4</b>	<b>71.9</b>	<b>13.5</b>
CN	<b>85.9</b>	<b>73.5</b>	<b>12.4</b>

**Table 4: Ablation study on Math23k.**

Models	Math23k
CN	<b>85.9</b>
CN w/o MT	85.2
CN w/o VIB	84.8
CN w/o MT+VIB	83.1
CN <sub>BERT</sub>	84.3
CN <sub>BERT</sub> w/o MT+VIB	82.4

**Table 5: Case study from Math23k. Demonstrates the ability of the model to generate multiple solutions.**

Problem Text:	A train leaves from place A at 7 o'clock and arrives at place B at 17 o'clock. The train travels 75 kilometers per hour. How many kilometers is the distance between the two places?
Ground-truth:	$(17 - 7) \times 75$
Source Network:	$(17 - 7) \times 75$
Collaborator Network:	$17 \times 75 - 7 \times 75$

the pre-training model, we also replaced the RoBERTa encoder with BERT [5] (denoted as CN<sub>BERT</sub>). We summarize the results in Table 4. Both the VIB strategy and mutual learning contribute greatly to the performance of *ESIB*.

### 4.3 Case Study

As an intuitive way to show the performance of *ESIB*, we randomly choose one problem from Math23k and show its solution expression generated by our model. As shown in Table 5, we observe that our model can produce two different but equivalent solutions. In addition, our model can correctly solve the problems in Table 1, which proves that our model does effectively reduce spurious correlations.

## 5 CONCLUSION

In this paper, we proposed an expression syntax information bottleneck method for MWP with mutual learning, which discards redundant features containing spurious correlations. Furthermore, we designed a self-distillation loss to encourage the model to rely more on the syntax information in the latent space. Extensive experiments on two benchmark MWP datasets demonstrated the effectiveness of our model.

## ACKNOWLEDGMENTS

This work was partially supported by Natural Science Foundation of Guangdong Province of China (No. 2021A1515011905), Shenzhen Science and Technology Innovation Program (Grant No. KQTD20190929172835662), Shenzhen Basic Research Foundation (No. JCYJ20210324115614039 and No. JCYJ20200109113441941)

## REFERENCES

- [1] Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. 2016. Deep variational information bottleneck. *arXiv preprint arXiv:1612.00410* (2016).
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [3] Daniel G Bobrow. 1964. Natural language input for a computer problem solving system. (1964).
- [4] Ting-Rui Chiang and Yun-Nung Chen. 2018. Semantically-aligned equation generation for solving and reasoning math word problems. *arXiv preprint arXiv:1811.00720* (2018).
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [6] Marco Federici, Anjan Dutta, Patrick Forré, Nate Kushman, and Zeynep Akata. 2020. Learning robust representations via multi-view information bottleneck. *arXiv preprint arXiv:2002.07017* (2020).
- [7] Ian Fischer. 2020. The conditional entropy bottleneck. *Entropy* 22, 9 (2020), 999.
- [8] Zhenya Huang, Qi Liu, Weibo Gao, Jinze Wu, Yu Yin, Hao Wang, and Enhong Chen. 2020. Neural mathematical solver with enhanced formula structure. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1729–1732.
- [9] Vivek Kumar, Rishabh Maheshwary, and Vikram Pudi. 2021. Adversarial Examples for Evaluating Math Word Problem Solvers. *arXiv preprint arXiv:2109.05925* (2021).
- [10] Zhenwen Liang and Xiangliang Zhang. 2021. Solving Math Word Problems with Teacher Supervision. *IJCAI* (2021).
- [11] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [12] Ilya Loshchilov and Frank Hutter. 2018. Fixing weight decay regularization in adam. (2018).
- [13] Sudipto Mukherjee, Himanshu Asnani, and Sreeram Kannan. 2020. CCMI: Classifier based conditional mutual information estimation. In *Uncertainty in artificial intelligence*. PMLR, 1083–1093.
- [14] Yu Pan, Zeyong Su, Ao Liu, Jingquan Wang, Nannan Li, and Zenglin Xu. 2022. A Unified Weight Initialization Paradigm for Tensorial Convolutional Neural Networks. In *ICML (Proceedings of Machine Learning Research, Vol. 162)*. PMLR, 17238–17257.

- [15] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP Models really able to Solve Simple Math Word Problems? *arXiv preprint arXiv:2103.07191* (2021).
- [16] Yibin Shen and Cheqing Jin. 2020. Solving math word problems with multi-encoders and multi-decoders. In *Proceedings of the 28th International Conference on Computational Linguistics*. 2924–2934.
- [17] Naftali Tishby, Fernando C Pereira, and William Bialek. 2000. The information bottleneck method. *arXiv preprint physics/0004057* (2000).
- [18] Lei Wang, Yan Wang, Deng Cai, Dongxiang Zhang, and Xiaojiang Liu. 2018. Translating a math word problem to an expression tree. *arXiv preprint arXiv:1811.05632* (2018).
- [19] Lei Wang, Dongxiang Zhang, Lianli Gao, Jingkuan Song, Long Guo, and Heng Tao Shen. 2018. Mathdqn: Solving arithmetic word problems via deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- [20] 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 *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 7144–7151.
- [21] Maolin Wang, Yu Pan, Xiangli Yang, Guangxi Li, and Zenglin Xu. 2023. Tensor Networks Meet Neural Networks: A Survey. *CoRR abs/2302.09019* (2023).
- [22] Yan Wang, Xiaojiang Liu, and Shuming Shi. 2017. Deep neural solver for math word problems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 845–854.
- [23] Qinzhuo Wu, Qi Zhang, Zhongyu Wei, and Xuan-Jing Huang. 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 (Volume 1: Long Papers)*. 5859–5869.
- [24] Zhipeng Xie and Shichao Sun. 2019. A Goal-Driven Tree-Structured Neural Model for Math Word Problems. In *IJCAI*. 5299–5305.
- [25] Jiong Xiong, Zixuan Li, Chuanyang Zheng, Zhijiang Guo, Yichun Yin, Enze Xie, Zhicheng Yang, Qingxing Cao, Haiming Wang, Xiongwei Han, et al. 2023. DQ-LORE: DUAL QUERIES WITH LOW RANK APPROXIMATION RE-RANKING FOR IN-CONTEXT LEARNING. *arXiv preprint arXiv:2310.02954* (2023).
- [26] Richard Zanibbi, Behrooz Mansouri, Anurag Agarwal, and Douglas W Oard. 2021. ARQMath: a new benchmark for math-aware CQA and math formula retrieval. In *ACM SIGIR Forum*, Vol. 54. ACM New York, NY, USA, 1–9.
- [27] Jipeng Zhang, Ka Wei LEE, Ee-Peng Lim, Wei Qin, Lei Wang, Jie Shao, Qianru Sun, et al. 2020. Teacher-student networks with multiple decoders for solving math word problem. (2020).
- [28] Jipeng Zhang, Lei Wang, Roy Ka-Wei Lee, Yi Bin, Yan Wang, Jie Shao, and Ee-Peng Lim. 2020. Graph-to-tree learning for solving math word problems. Association for Computational Linguistics.
- [29] Ying Zhang, Tao Xiang, Timothy M Hospedales, and Huchuan Lu. 2018. Deep mutual learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4320–4328.
- [30] Wei Zhao, Mingyue Shang, Yang Liu, Liang Wang, and Jingming Liu. 2020. Ape210k: A large-scale and template-rich dataset of math word problems. *arXiv preprint arXiv:2009.11506* (2020).