

# 基于指令微调的数学推理任务探索

Huazheng Zeng  
School of Computer Science and Technology  
Fudan University  
220 Handan Rd, Shanghai, 200433  
22302010022@m.fudan.edu

Yang Wang  
School of Computer Science and Technology  
Fudan University  
220 Handan Rd, Shanghai, 200433  
223020100@m.fudan.edu

Yuxuan Xie  
School of Computer Science and Technology  
Fudan University  
220 Handan Rd, Shanghai, 200433  
223020100@m.fudan.edu

## Abstract

1 数学推理是评估人类智力基本认知能力的基石。近年来，针对自动解决数  
2 学问题的大型语言模型 LLMs 得到了显著发展。我们基于 Qwen2.5-0.5B-  
3 instruction，探索了全参数 SFT，全参数 LoRA 指令微调的效果，同时，通  
4 过数据集增强，提高了模型的性能；我们还在解码阶段加入多轮投票，得  
5 到了相关的测试结果。

## 6 1 Introduction

7 数学推理是评估人类智力的核心任务之一。随着大规模语言模型 (LLMs) 的不断发展，基  
8 于自然语言处理的数学推理任务已成为一个重要的研究领域。尽管传统的数学推理方法依  
9 赖于符号计算和规则推导，近年来，深度学习，尤其是大型语言模型，已经在自动推理方面  
10 取得了显著的进展。

11 本研究基于 Qwen2.5-0.5B-instruction 模型，探讨了指令微调 (instruction fine-tuning) 方  
12 法在数学推理任务中的应用，重点研究了全参数微调 (SFT) 和 LoRA 微调两种技术的效  
13 果。此外，我们还结合了数据增强策略，通过 DeepSeek API 扩展训练数据集，以提升模型  
14 在复杂数学推理任务中的表现。为进一步增强模型的推理能力，我们在解码阶段引入了多轮  
15 投票机制，旨在通过集成多次推理结果来提高推理的稳定性和准确性。

16 本文的主要贡献包括：提出了 Qwen2.5 模型在数学推理任务中的应用，探索了不同微调技  
17 术的效果，并评估了数据增强和投票机制在该任务中的有效性。

## 18 2 Related Work

### 19 2.1 数据增强与一致性训练

20 数据增强在提升机器学习模型泛化能力方面起到了关键作用，尤其是在低资源场景下。Asai  
21 和 Hajishirzi (2020) 提出了基于逻辑的数据增强与正则化 (Logic-Guided Data Augmentation  
22 and Regularization, LGDA)，通过结合逻辑规则与数据增强，提高了问答 (QA) 任务的准  
23 确性和一致性。他们的方法通过对称性和传递性等逻辑规则生成新的训练样本，并通过正则  
24 化约束确保模型预测的一致性，从而增强了模型的推理能力。在 WIQA 和 QuaRel 等基准  
25 数据集上，LGDA 显著减少了预测不一致性，并提升了模型性能。

26 LGDA 的逻辑驱动数据增强策略启发我们在数据增强中探索隐藏的逻辑结构，而其一致性  
27 正则化的创新进一步强调了确保预测连贯性的重要性。这些思想也可以扩展到其他需要复  
28 杂推理能力的领域中。

### 29 2.2 数据增强在半监督学习中的作用

30 Xie 等人 (2020) 提出了无监督数据增强 (Unsupervised Data Augmentation, UDA)，展示  
31 了高质量数据增强技术在半监督学习一致性训练框架中的显著效果。他们的方法引入了诸  
32 如 RandAugment 和反向翻译等先进的数据增强技术，在有限标注数据的条件下显著提升  
33 了模型性能。UDA 在多种语言和视觉任务上实现了最新的性能水平，强调了数据增强的多  
34 样性和自然性在提升模型鲁棒性方面的重要性。

35 UDA 强调了增强数据多样性在提高模型一致性中的关键作用。这一原则启发我们使用多样  
36 化和高质量的数据增强策略，以提升模型在复杂任务（例如数学推理或其他深度学习领域）  
37 中的表现。

### 38 2.3 启发

39 LGDA 和 UDA 在数据增强与一致性训练方面的研究为我们的工作设计提供了重要的启发。  
40 LGDA 通过逻辑规则生成增强样本的策略激发了我们在领域特定任务中融入逻辑结构的潜  
41 力，而 UDA 在利用高质量数据增强技术提升模型鲁棒性方面的成功，则进一步证明了数据  
42 增强在提高模型能力中的不可或缺性。这些方法共同表明，即使在有限标注数据的条件下，  
43 通过数据增强与一致性训练，模型的推理能力和泛化性能也可以显著提高。

## 44 3 Method

45 在本研究中，我们采用了多种方法来提升 Qwen 系列模型在数学推理任务上的表现。我们  
46 的方法结合了以下主要技术：

- 47 • 1 全参数微调，增强模型在特定任务上的适应能力；
- 48 • 2 LoRA 微调，优化计算资源效率的同时能取得较好的性能；
- 49 • 3 使用 DeepSeek 在原有数据集上生成回答以对其进行扩充；
- 50 • 4 多数投票推理机制，以期提升推理过程的稳健性；
- 51 • 5 数据增强策略，进一步提升模型的泛化能力。

### 52 3.1 全参数微调

全参数微调对模型的全部参数进行优化，以提升模型在特定任务（如数学推理任务）上的适  
应能力：

$$\mathcal{L}_{full} = \frac{1}{N} \sum_{i=1}^N \ell(f(X_i; \Theta), y_i)$$

53 其中， $X_i$  表示输入问题， $y_i$  表示期望输出， $\Theta$  为模型的全部参数。

### 54 3.2 LoRA 微调

为减少全参数微调带来的计算成本，同时降低其灾难性遗忘风险，以提高模型对新任务的  
迁移能力，我们采用了 LoRA 技术进行微调。该技术在微调过程中固定大部分预训练参数，

仅更新小规模适配器矩阵：

$$\Theta' = \Theta + A \cdot B$$

55 假设  $\Theta \in \mathbb{R}^{d \times k}$ ，则  $A \in \mathbb{R}^{d \times r}$ 、 $B \in \mathbb{R}^{r \times k}$ ，且  $r \ll \min(d, k)$ 。

### 56 3.3 DeepSeek 生成训练数据

57 DeepSeek 是一个较为强大的大语言模型，其在数学推理方面的能力也较为出色。故我们尝试  
58 调用 DeepSeek api 对 GSM8K 和 MATH 训练集生成回答以扩充训练数据。由于 DeepSeek  
59 生成的回答能够包含较为详细的解答步骤，故我们期望通过这种类似蒸馏的过程来指导  
60 Qwen2.5B 提升数学推理方面的能力。

### 61 3.4 多数投票推理机制

62 在推理阶段，我们通过多数投票机制以期增强模型推理的鲁班性。假设模型的正确率为  $p$ ，  
63 在多数投票推理阶段生成  $r$  次回答，则期望的正确回答数量为  $p \times r$ 。若考虑最佳情况，其  
64 余错误回答均不一致，那么希望从该机制中获得正确回答，则需要  $p \times r \geq 2$ ，即  $p \geq \frac{2}{r}$ 。但  
65 在实际测试中，该方法反而降低了模型的准确性，这可能是由于模型生成的回答与最佳情况  
66 的假设不符，生成了相同的错误答案。

### 67 3.5 数据增强

68 由于数据在模型微调中占有较为重要的作用，所以我们采用了《MuggleMath: Assessing  
69 the Impact of Query and Response Augmentation on Math Reasoning》一文中提供的  
70 AugGSM8K 和 AugMATH 数据集对模型进行微调，并对其进行过滤以保留高质量训练数  
71 据，以此让模型更好地学习数学推理能力，并取得了较好的效果。

## 72 4 Evaluation

### 73 4.1 基于 SFT 与 LoRA 的微调对比

74 在本部分中，我们比较了全参数微调（SFT）和 LoRA 微调的效果。我们基于 Qwen2-0.5B  
75 模型对 SFT 和 LoRA 两种微调方法进行了测试，并对比了显存占用、训练时间和训练效果。  
76 实验结果如下表所示：

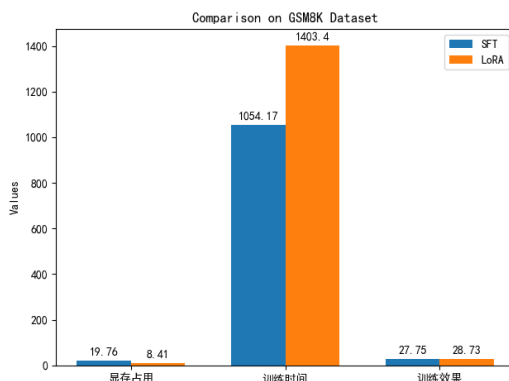


Figure 1: 基于 GSM8K 数据集的 SFT 与 LoRA 微调对比

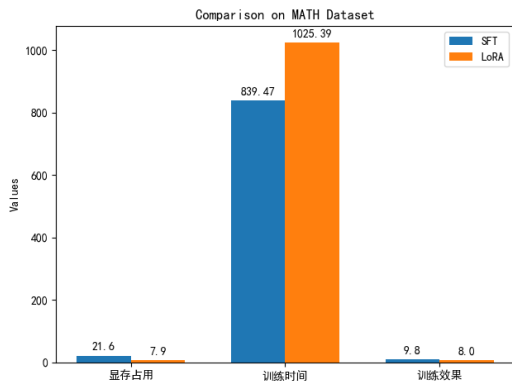


Figure 2: 基于 MATH 数据集的 SFT 与 LoRA 微调对比

分析：全参数微调 (SFT) 在训练大规模任务时效果较好，但由于需要更新模型的所有参数，其显存占用较高且训练时间较长。相比之下，LoRA 微调仅更新少部分适配器矩阵，显著降低了显存占用和训练时间，同时在小规模任务中表现优异，且能够较好地适应不同任务。

#### 4.2 基于增强数据集训练与正常训练的微调对比

在这部分实验中，我们对比了基于增强数据集训练和正常训练的微调效果。我们通过 DeepSeek 生成的扩充数据集以及 AugGSM8K 和 AugMATH 数据集进行了微调，并评估了训练效果。

##### 4.2.1 基于 MuggleMath 的数据增强训练

我们基于 MuggleMath 数据，通过正则匹配消去重复数据，并且按一定概率采样，最后提取出两个个 15k 条的数据集 AugGSM8K 和 AugMATH 数据集，用于对比微调效果：

Table 1: 基于 AugGSM8K 训练与正常训练的微调对比

指标	正常训练	增强数据集训练
训练时间	1054.17s	2438.16s
训练效果	27.75%	42.99%

Table 2: 基于 AugMATH 增强数据集训练与正常训练的微调对比

指标	正常训练	增强数据集训练
训练时间	839.47s	2099.35s
训练效果	9.80%	12.00%

#### 分析：

结果显示，我们已经超越了基于原始数据集 baseline 的训练效果。增强数据集的使用能够显著提升模型在数据多样性和鲁棒性上的表现。尽管增强数据集训练的时间较长，但其在训练效果上超越了常规训练，尤其是在数学推理等复杂任务中，增强数据集对提升模型的泛化能力具有显著作用。

同时，在构造数据集的过程中，我们发现包含更多符号推导的数据集更能够增强模型的推理能力，相对于使用自然语言推导而言。使用自然语言推导的增强数据集，反而会使得模型推理性能减弱，同时最后很难去匹配正确答案。

##### 4.2.2 DeepSeek 生成的数据集

我们使用 DeepSeek 生成的回答来扩充 GSM8K 和 MATH 训练集。DeepSeek 能够生成带有详细解答步骤的回答，这种类似蒸馏的过程帮助 Qwen2.5B 模型提升了数学推理方面的能力。生成的数据集包含 15k 条数据，但我们没有做更多处理。

Table 3: deepseek-基于增强数据集训练与正常训练的微调比对

指标	正常训练	增强数据集训练
训练时间	1054.17s	1393.15s
训练效果	27.75%	25.55%

99 **分析：**  
100 这个结果表明，简单地添加更多的数据并不一定能够提升模型的表现，尤其是当这些数据是  
101 通过某种生成方式获得的时候。在机器学习和深度学习领域，数据的质量往往比数量更为重  
102 要。增强数据集需要有技巧地构造，以确保新加入的数据能够有效地补充原有数据集的不  
103 足，提供额外的信息或多样性，而不是简单地重复或引入噪声。

104 **4.3 更大模型验证**

105 为了验证我们的数据集的有效性和广泛适用性，我们使用 3B 模型进行了验证

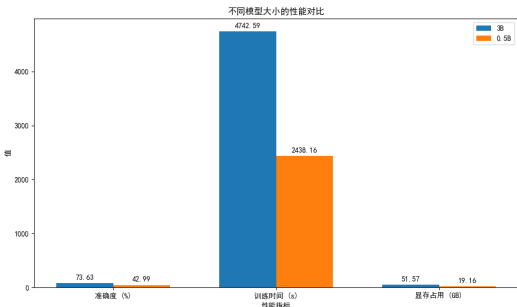


Figure 3: 3B 模型与 0.5B 模型的增强数据集训练对比

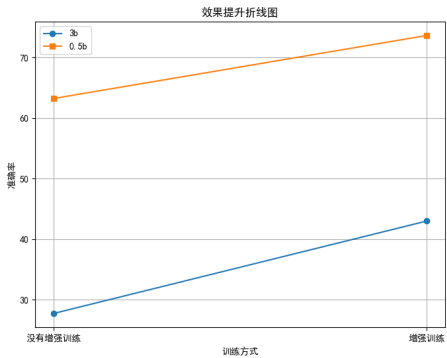


Figure 4: 效果曲线图

106 **分析：**  
107 通过更大模型的验证，证明了我们的增强数据集的有效性和广泛适应性。但更大的模型提高  
108 更加有限，还是需要力大砖飞

109 **4.4 多轮投票机制的比对**

Table 4: 多轮投票机制的比对

指标	无投票	一轮投票	二轮投票	三轮投票
训练效果	28.73%	29.69%	0%	0%

在多轮投票机制中，我们期望通过生成多个回答并进行投票来提升推理的鲁棒性。然而，实验结果表明，多轮投票机制并未显著提升模型的准确性，甚至在某些情况下导致训练效果下降。以下对这一现象进行详细分析：

### 1. 错误答案的重复性

多轮投票的基本假设是生成多个不同的答案，通过多数投票排除错误。然而，当模型生成的回答存在高度相似的错误时，投票机制无法有效纠正错误。重复的错误答案在投票中占据多数，导致最终结果仍然错误，失去了投票机制的预期优势。

### 2. 模型固有偏差的影响

由于模型训练数据或参数设置的固有偏差，生成的答案容易呈现出一致的错误模式。这种情况下，多轮生成并未引入足够的多样性，导致相似的错误答案被多次投票确认，进一步加剧了错误结果的稳定性。

### 3. 生成结果的多样性不足

多轮投票的有效性依赖于生成结果的多样性，即多个回答之间需存在一定的逻辑或内容差异。然而，实验中模型生成的回答往往高度一致，缺乏足够的多样性，使得投票机制无法剔除错误答案。此外，模型对问题的理解局限性也限制了结果的生成范围。

### 4. 训练资源与优化设计限制

从实验结果来看，二轮投票和三轮投票的训练效果均为 0，可能反映出投票机制设计与模型优化之间存在不匹配。训练过程中资源分配不均或多轮生成策略未经过有效的优化，导致投票机制无法发挥应有的作用。

## 5 Conclusion

本研究探索了多种方法提升 Qwen 系列模型在数学推理任务上的表现。实验结果表明，基于全参数微调 (SFT) 和 LoRA 微调的模型在数学推理任务中表现出色，尤其是在特定任务下，通过优化计算资源能够显著提升模型的推理能力。同时，数据增强策略有效提升了模型在复杂任务中的泛化能力，尤其是在基于 DeepSeek 生成的扩充数据集进行训练后，模型的数学推理能力得到了明显提升。

然而，实验结果也表明，多轮投票机制在提高推理精度方面的效果并不显著，可能是由于生成的多个答案存在重复的错误。因此，在模型推理过程中的策略选择仍需进一步优化。此外，随着模型规模的增加，显存占用和训练时间也呈现上升趋势，如何在性能和资源消耗之间找到最佳平衡仍然是一个值得关注的研究方向。

未来的研究可以从以下几个方面进行拓展：进一步优化微调策略，探索不同类型的数据增强方法，提升多轮投票机制的鲁棒性，探索如何结合符号推理与深度学习模型来增强数学推理能力，以及多步骤的微调。

## 6 References

- [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.
- [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural SIMulation System*. New York: TELOS/Springer-Verlag.
- [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* 15(7):5249–5262.