

# MetaMath: Bootstrap Your Own Mathematical Questions for Large Language Models

Longhui Yu<sup>1,\*</sup>, Weisen Jiang<sup>2,3,\*</sup>, Han Shi<sup>4</sup>, Jincheng Yu<sup>3,4</sup>,  
Zhengying Liu<sup>4</sup>, Yu Zhang<sup>2</sup>, James T. Kwok<sup>3</sup>, Zhenguo Li<sup>4</sup>,  
Adrian Weller<sup>1,5</sup>, Weiyang Liu<sup>1,6</sup>

<sup>1</sup>University of Cambridge <sup>2</sup>SUSTech

<sup>3</sup>HKUST <sup>4</sup>Huawei Noah's Ark Lab

<sup>5</sup>The Alan Turing Institute <sup>6</sup>Max Planck Institute

## Methods

### Answer Augmentation (AnsAug)

**Question:** James buys 5 packs of beef that are 4 pounds each. The price of beef is \$5.50 per pound. How much did he pay?

**Answer:** (sample answers from GPT-3.5-Turbo)

### Rephrasing Question

**Question:** What is the total amount that James paid when he purchased 5 packs of beef, each weighing 4 pounds, at a price of \$5.50 per pound?

**Answer:** Each pack of beef weighs 4 pounds, so 5 packs weigh  $4 * 5 = 20$  pounds in total. The price per pound of beef is \$5.50, so the total cost for 20 pounds is  $20 * \$5.50 = \$110$ . ... The answer is: 110.

### Backward Question by Self-Verification (needs rewriting) (EMNLP 2023)

**Question:** James buys  $x$  packs of beef that are 4 pounds each. The price of beef is \$5.50 per pound. **He paid 110. What is the value of unknown variable  $x$ ?**

**Answer:** To solve this problem, we need to determine the value of  $x$ , which represents the number of packs of beef that James bought. Each pack of beef weighs 4 pounds and ... The value of  $x$  is 5.

### Backward Question by FOBAR (Preprint 2023)

**Question:** James buys  $x$  packs of beef that are 4 pounds each. The price of beef is \$5.50 per pound. How much did he pay? **If we know the answer to the above question is 110, what is the value of unknown variable  $x$ ?**

**Answer:** James buys  $x$  packs of beef that are 4 pounds each, so he buys a total of  $4x$  pounds of beef. The price of beef is \$5.50 per pound, so the total cost of the beef is  $5.50 * 4x = 22x$ . ... The value of  $x$  is 5.

**Datasets:** 395K augmented samples by 4 methods on original GSM8K & MATH training data

| Dataset          | AnsAug | Rephrasing | SV  | FOBAR | Overall |
|------------------|--------|------------|-----|-------|---------|
| MetaMathQA-GSM8K | 80K    | 80K        | 40K | 40K   | 240K    |
| MetaMathQA-MATH  | 75K    | 50K        | 15K | 15K   | 155K    |
| MetaMathQA       | 155K   | 130K       | 55K | 55K   | 395K    |

**Ablation:** data augmented by 4 methods all have performance gains.

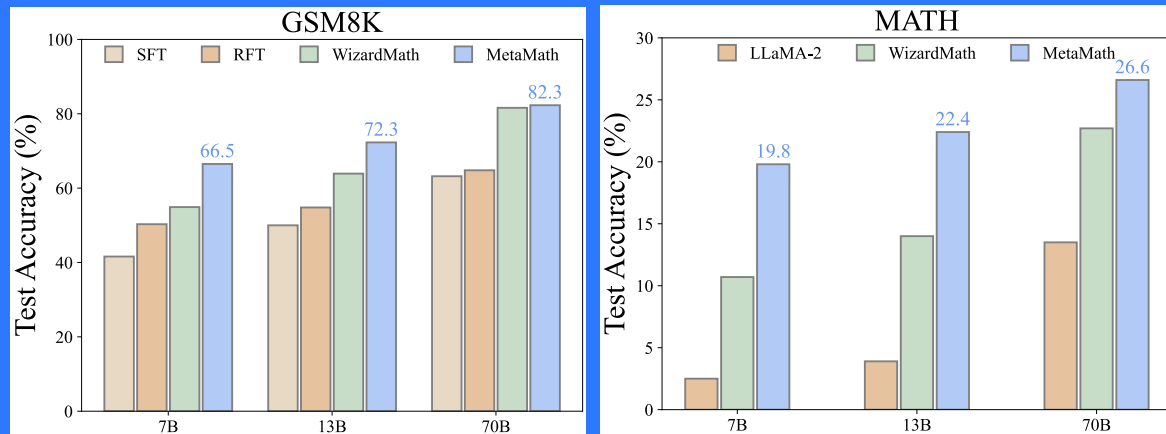
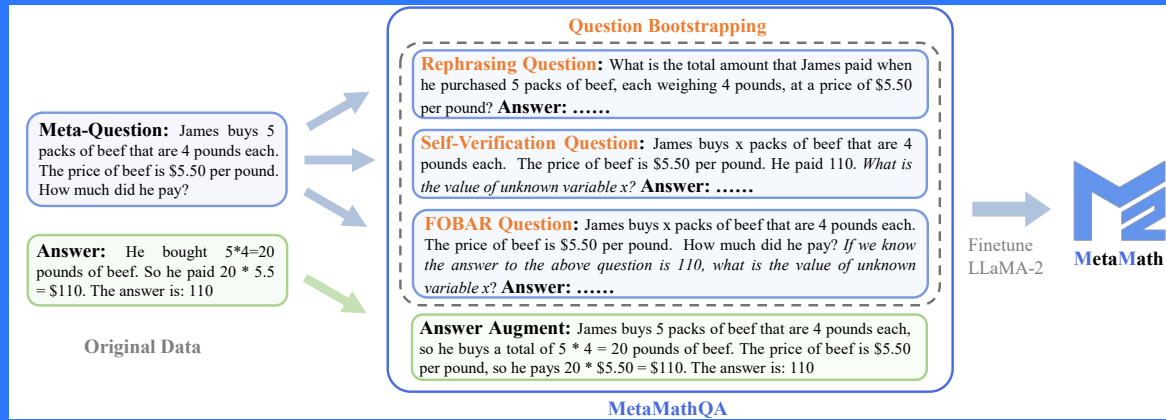
| Method   | GSM8K  |      |    |       |             |      | MATH   |      |    |       |             |             |
|----------|--------|------|----|-------|-------------|------|--------|------|----|-------|-------------|-------------|
|          | AnsAug | Rep. | SV | FOBAR | GSM8K       | MATH | AnsAug | Rep. | SV | FOBAR | GSM8K       | MATH        |
| SFT      | ✗      | ✗    | ✗  | ✗     | 41.6        | 3.0  | ✗      | ✗    | ✗  | ✗     | 13.8        | 4.7         |
| MetaMath | ✓      | ✗    | ✗  | ✗     | 59.6        | 4.4  | ✓      | ✗    | ✗  | ✗     | 28.4        | 12.9        |
|          | ✗      | ✓    | ✗  | ✗     | 59.7        | 4.4  | ✗      | ✓    | ✗  | ✗     | 30.4        | 12.4        |
|          | ✓      | ✓    | ✗  | ✗     | 60.6        | 4.4  | ✓      | ✓    | ✗  | ✗     | 29.1        | 15.3        |
|          | ✓      | ✓    | ✓  | ✓     | <b>64.4</b> | 5.7  | ✓      | ✓    | ✓  | ✓     | <b>34.6</b> | <b>17.7</b> |

## Spotlight

We improve LLM's Math reasoning ability in both forward reasoning & backward reasoning



Scan to  
code / data /  
checkpoints



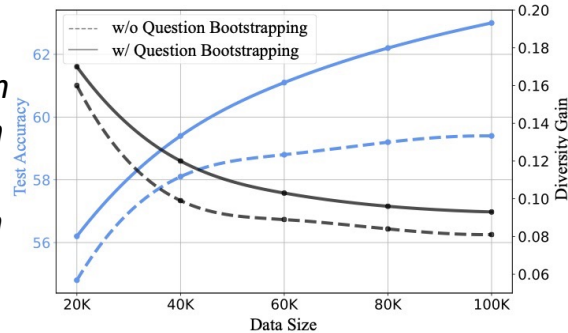
UNIVERSITY OF  
CAMBRIDGE



MAX-PLANCK-GESELLSCHAFT

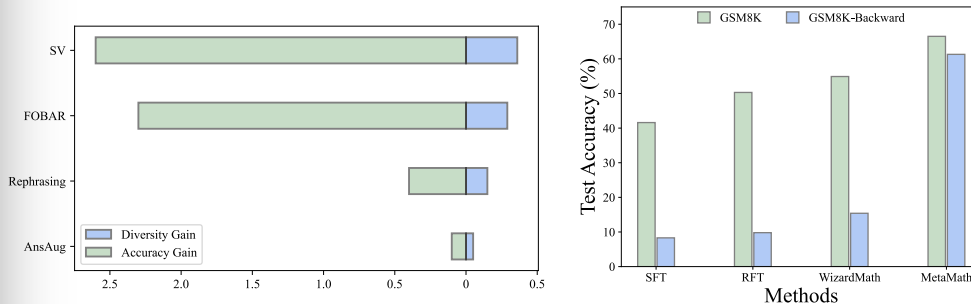
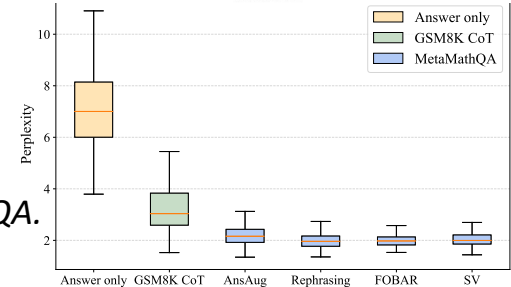
### Diversity & Accuracy:

Naïve data augmentation suffers a quick saturation in accuracy. Thanks to high diversity, MetaMath alleviates saturation.



### Lower perplexity of MetaMathQA:

Pretrained models (e.g. LLaMA-2) have lower perplexity on MetaMathQA.



**Left:** SV & FOBAR bring higher diversity, resulting in higher performance gains. **Right:** MetaMath achieves better backward reasoning ability than existing methods.

### MetaMath on OOD tasks & stronger models:

#### Performance on DROP dataset Performance on Llemma & Mistral

Different from GSM8K & MATH, On stronger models such as Llemma & Mistral, Metamath also boosts performance.

|            | #Params | Accuracy (Exact Match) |
|------------|---------|------------------------|
| SFT        | 7B      | 25.8                   |
| RFT        | 7B      | 26.7                   |
| WizardMath | 7B      | 31.5                   |
| MetaMath   | 7B      | <b>37.1</b>            |
| WizardMath | 13B     | 46.4                   |
| MetaMath   | 13B     | <b>49.5</b>            |
| WizardMath | 70B     | 63.1                   |
| MetaMath   | 70B     | <b>72.3</b>            |

|             | MetaMathQA | GSM8K       | MATH        |
|-------------|------------|-------------|-------------|
| LLaMA-2-7B  | ✗          | 14.6        | 2.5         |
|             | ✓          | <b>66.5</b> | <b>19.8</b> |
| LLaMA-2-13B | ✗          | 28.7        | 3.9         |
|             | ✓          | <b>72.3</b> | <b>22.4</b> |
| Llemma-7B   | ✗          | 36.4        | 18.0        |
|             | ✓          | <b>69.2</b> | <b>30.0</b> |
| Mistral-7B  | ✗          | 52.2        | 13.1        |
|             | ✓          | <b>77.7</b> | <b>28.2</b> |