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

Longhui Yu[1,*], Weisen Jiang[2,3,*], Han Shi[4], Jincheng Yu[3,4], Zhengying Liu[4], Yu Zhang[2], James T. Kwok[3], Zhenguo Li[4], Adrian Weller[1,5], Weiyang Liu[1,6]

[1]University of Cambridge [2]SUSTech
[3]HKUST [4]Huawei Noah's Ark Lab
[5]The Alan Turing Institute [6]Max Planck Institute

## Methods

**Answer Augmentation (AnsAug)**

**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:** (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.

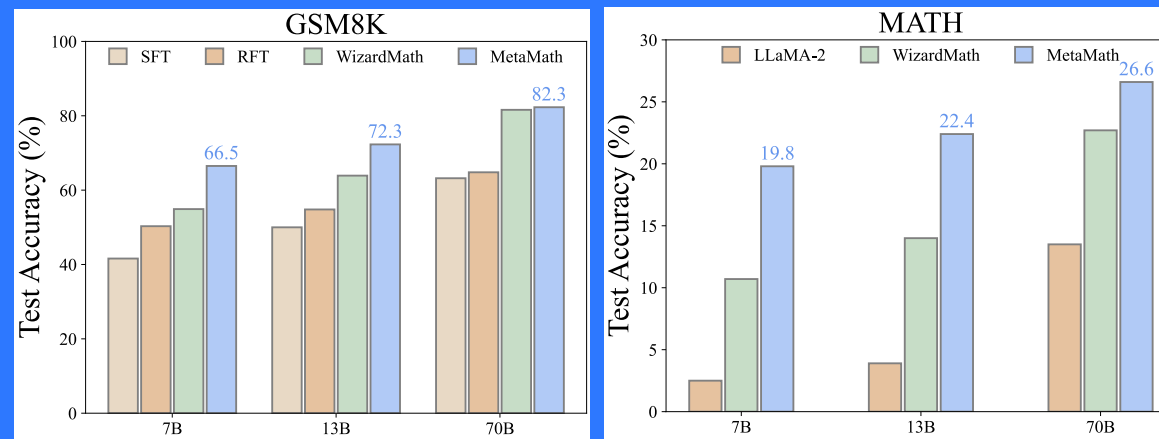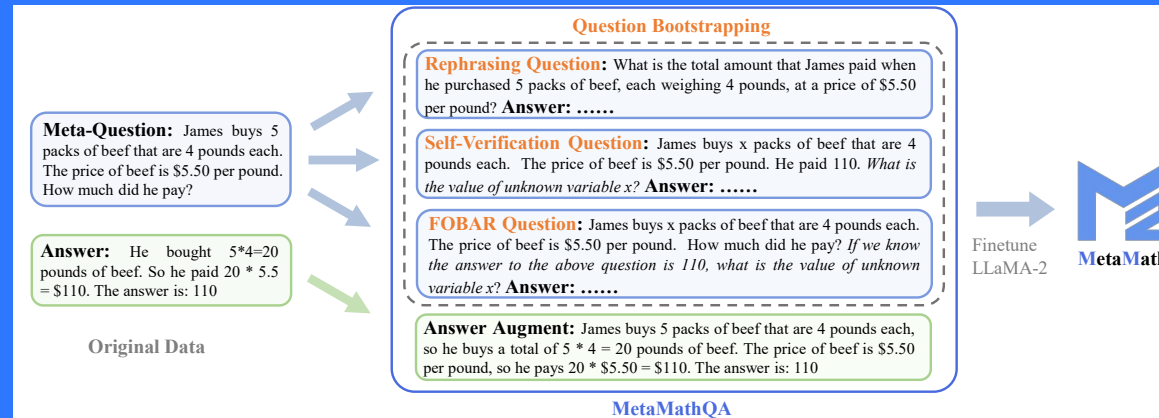**Datasets:** 395K augmented samples by 4 methods on original GSM8K & MATH training data

| 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.

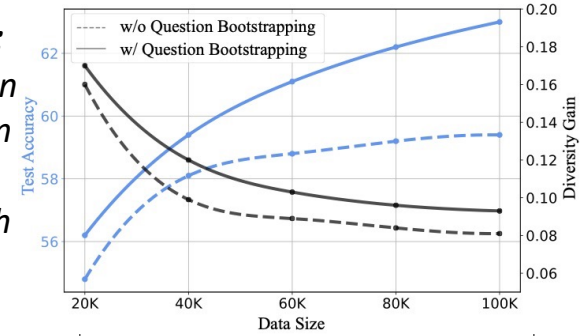| Method | GSM8K | | | | | | MATH | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | AnsAug | Rep. | SV | FOBAR | GSM8K | MATH | AnsAug | Rep. | SV | FOBAR | GSM8K | MATH |
| SFT | ✗ | ✗ | ✗ | ✗ | 41.6 | 3.0 | ✗ | ✗ | ✗ | ✗ | 13.8 | 4.7 |
| MetaMath | ✓ | ✗ | ✗ | ✗ | 59.6 | 4.4 | ✓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

# Spotlight

We improve LLM's **Math** reasoning ability in both

**forward reasoning &**

**backward reasoning**



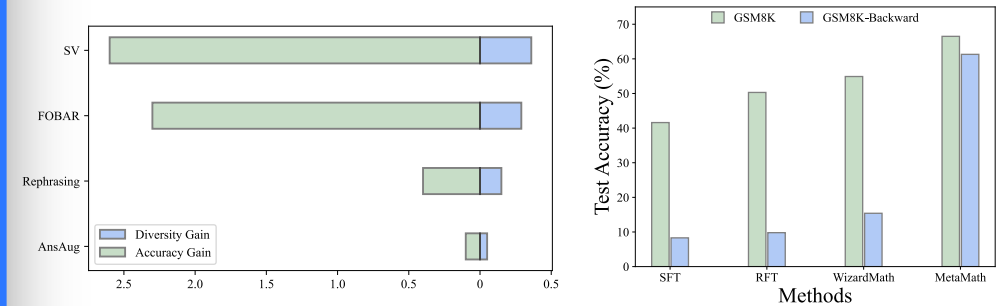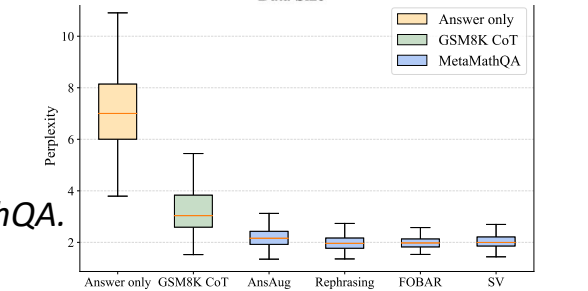Scan to **code / data / checkpoints**





**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**
*Different from GSM8K & MATH, questions in DROP have longer reasoning context. MetaMath performs better than baselines.*

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

**Performance on Llemma & Mistral**
*On stronger models such as Llemma & Mistral, Metamath also boosts performance.*

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