



香港城市大學  
City University of Hong Kong



جامعة محمد بن زايد  
للذكاء الاصطناعي  
MOHAMED BIN ZAYED UNIVERSITY  
OF ARTIFICIAL INTELLIGENCE



**ICLR**



# MUSTARD: Mastering Uniform Synthesis of Theorem and Proof Data

Yinya Huang<sup>1,6</sup>, Xiaohan Lin<sup>2</sup>, Zhengying Liu<sup>3</sup>, Qingxing Cao<sup>2</sup>, Huajian Xin<sup>2</sup>, Haiming Wang<sup>2</sup>,  
Zhenguo Li<sup>3</sup>, Linqi Song<sup>1,6</sup>, Xiaodan Liang<sup>2,4,5</sup>

<sup>1</sup>City University of Hong Kong <sup>2</sup>Shenzhen Campus of Sun Yat-sen University

<sup>3</sup>Huawei Noah's Ark Lab <sup>4</sup>DarkMatter AI Research <sup>5</sup>MBZUAI <sup>6</sup>CityUSRI

Contact: [yinya.huang@hotmail.com](mailto:yinya.huang@hotmail.com)



# Intermediate Step Validation

- Correct intermediate steps are crucial for LLMs to perform complex reasoning.
- High-quality step-wise annotations are hard to obtain.
- Current data synthesis trades off between correctness and reusability.
- ✓ MUSTARD synthesizes **correct, scalable, and reusable** mathematical data by combining the advantages of LLMs in verbalization and formal theorem provers in rigorous data validation.

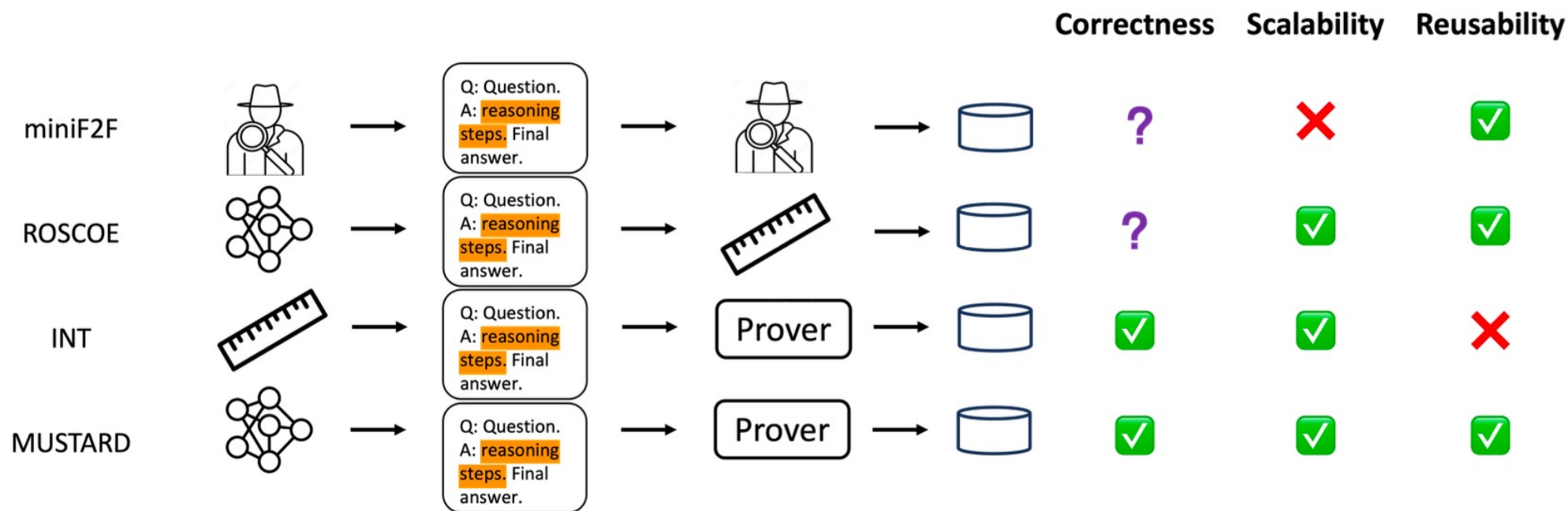
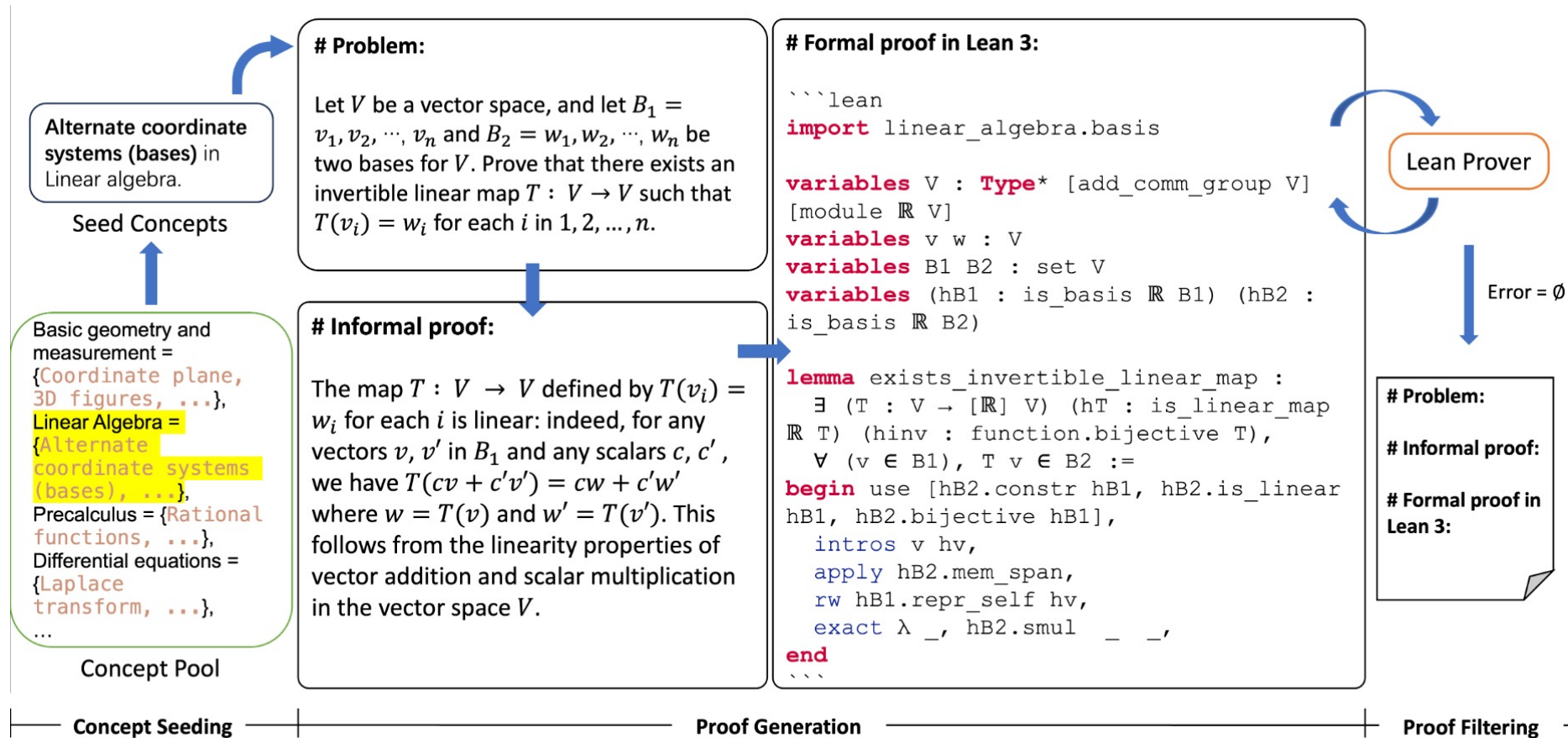


Figure 1: A comparison of methods of synthesizing and validating intermediate reasoning steps.



- Sampling math concepts from concept pool
- Prompting an LLM to generate problem, informal solution, and formal solution in Lean 3.
- Collecting data by Lean Prover validation.

**MUSTARDSAUCE** (MUSTARD resource)



- ✓ 5,866 MWP & ATP data points w/ {informal, formal} {statement, solution}.
- ✓ Proof length increases over educational levels.
- ✓ The most challenging problems require around 30 proof steps w/ 20 Lean tactics.

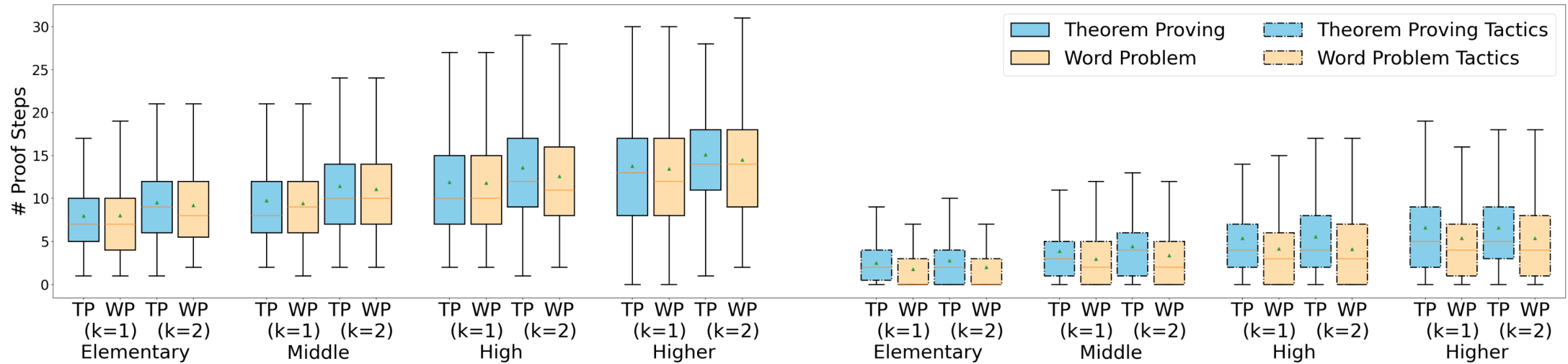


Figure 4: Distributions of formal proof lengths.



- We randomly select 200 data points from the generated data, 100 of which pass the Lean Prover (Group Valid) and 100 of which do not (Group Invalid).
- 6 inspection dimensions including factuality check and consistency check on informal/formal statement/proof.
- ✓ **(D1)** shows significance. → High quality math questions.
- ✓ **(D4)** and **(D6)** show significant differences between Group Valid & Invalid. → High quality autoformalization.

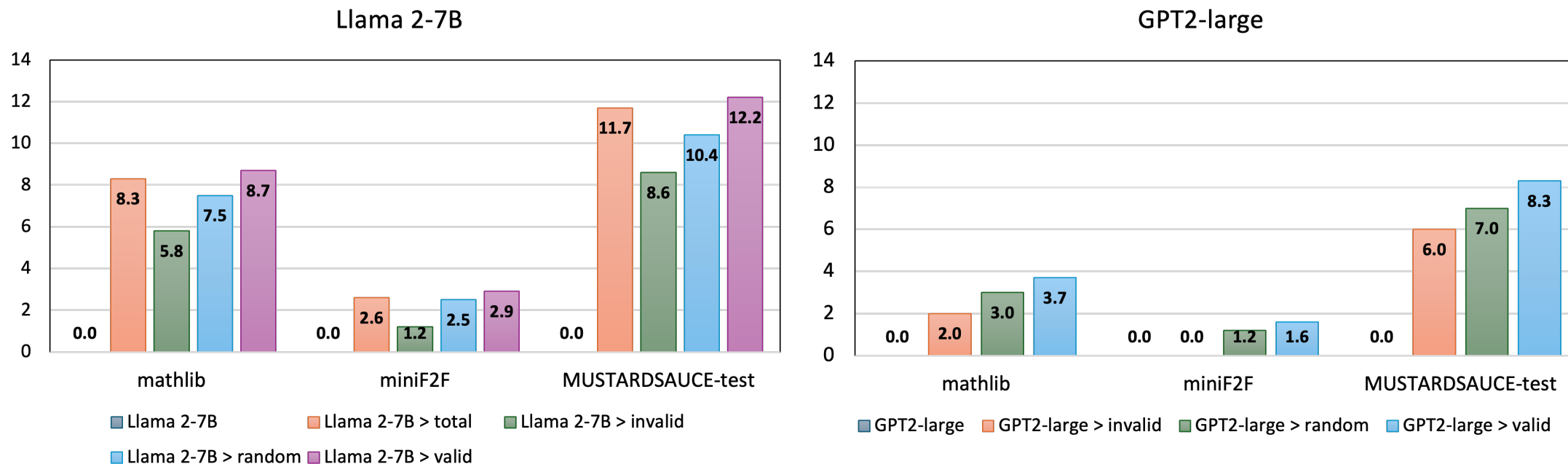
Table 3: Inspection dimensions and requirements in human evaluation. IS: Informal Statement. IP: Informal Proof. FS: Formal Statement. FP: Formal Proof. RT: Reasoning Type. Significant  $p < 0.005$  are marked with **bold**.

Inspection Dimension	Requirement	Valid	Invalid	$p$ -value
(D1) IS Correctness	<i>Whether the informal statement is factually correct.</i>	93.50	83.50	<b>0.00167</b>
(D2) IS Relevance	<i>Whether the informal statement is relevant to each seed concept.</i>	87.50	92.50	0.09604
(D3) RT Classification	<i>Whether the informal statement is of the required question type.</i>	67.00	68.50	0.74903
(D4) IP Correctness	<i>Whether the informal proof correctly solves the informal statement.</i>	88.50	73.50	<b>0.00012</b>
(D5) IS-FS Alignment	<i>Whether the informal statement and the formal statement describe the same problem and are aligned with each other.</i>	74.00	66.50	0.10138
(D6) IP-FP Alignment	<i>Whether the informal proof and the formal proof describe the same solution and have aligned proof steps.</i>	72.00	54.00	<b>0.00018</b>



- Using the generated MUSTARDSAUCE to fine-tune LMs and then validate on automated theorem proving (ATP).
- ✓ Average of **18.15%** relative performance gain.
- ✓ Llama 2-7B achieves **16.00%** gain on mathlib, **16.00%** gain on miniF2F, and **17.31%** gain on MUSTARDSAUCE-test.

**Fig.** Pass@1 results on automated theorem proving (ATP) tasks. > denotes a fine-tuning step. test: MUSTARDSAUCE-test. Note that the reported results on MUSTARDSAUCE-test are obtained by only fine-tuning on the MUSTARDSAUCE-valid training split.



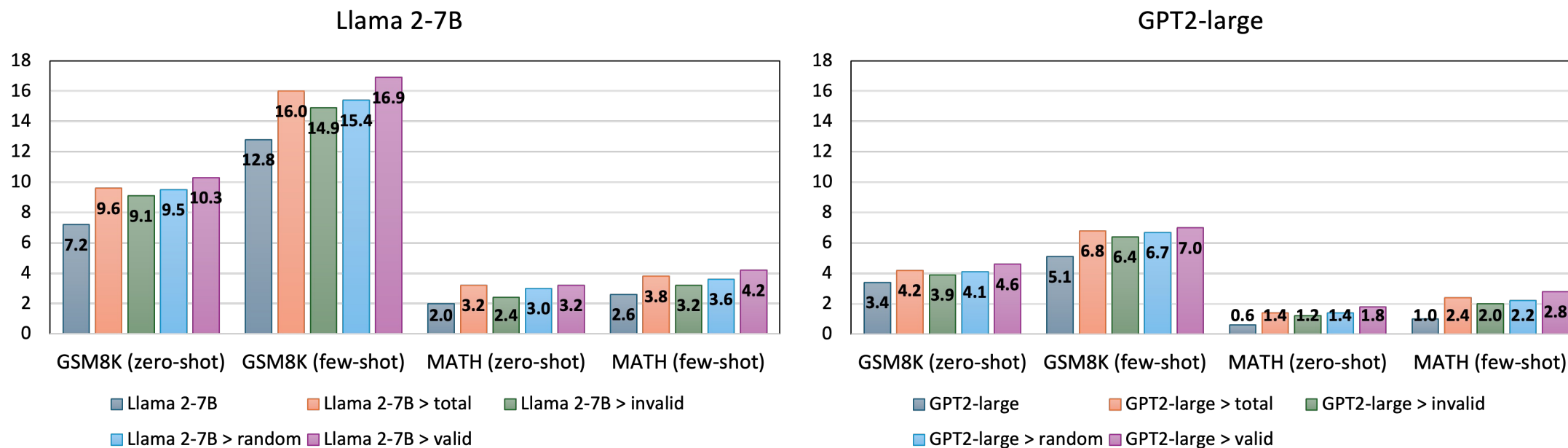


# Data Benefits: Math Word Problem



- Using the generated MUSTARDSAUCE to Fine-tune LMs and then validate on math word problems (MWP).
- ✓ Average of **11.01%** relative performance gain.
- ✓ GPT2-large achieves **28.57%** gain on MATH, **12.20%** gain on GSM8K, zero-shot setting.

Fig. Maj1@1 results on GSM8K (G) and MATH (M). Zero: Zero-shot. Few: Few-shot. > denotes a fine-tuning step.





## Poster: Halle B, Tue 7 May 4:30— 6:30 p.m. CEST

1. We propose a novel framework **MUSTARD** that can generate high-quality mathematical data (both informal and formal) with an interplay between generative language model and theorem prover assistants.
2. We release the **MUSTARDSAUCE**, which contains both math word problems and theorem-proving problems spanning elementary to higher educational levels. Each sample has corresponding informal and formal solutions.
3. We conduct extensive analysis and experiments on the generated data, demonstrating their quality, diversity, and effectiveness in improving language models' mathematical reasoning performance.

Contact: [yinya.huang@hotmail.com](mailto:yinya.huang@hotmail.com)



👉 Read our paper



👉 Try our code



👉 Take on the exciting challenges!

