# Programming every example:
# Lifting pre-training data quality like experts at scale

**Fan Zhou [1,3]*, Zengzhi Wang[1,3]*, Qian Liu[2], Junlong Li[1], Pengfei Liu[1,3]**
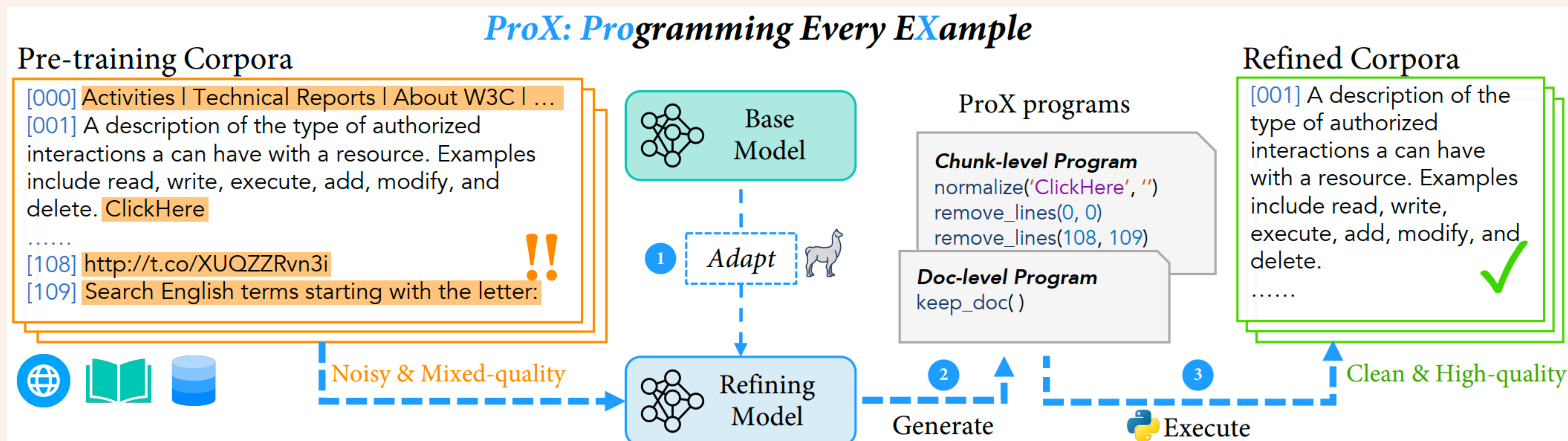
[1]Shanghai Jiao Tong University  [2]Sea AI Lab  [3]GAIR Lab

🫐 **TL;DR** Still worried about the potential noise and the low quality of your rule-cleaned pre-training corpora? **Try ProX!** In ProX, we use 0.3B LLMs to seamlessly refine your pre-training dataset, providing a clean start for your LLM training. ProX serves as a LLM-driven framework which generate **sample-wise**, **executable** cleaning programs to clean **each and every one** of your pre-training data samples.

- Our 1.7B model, trained on ProX corpus with 50B tokens training, performs on par with TinyLlama-1.1B which is trained on 3T tokens.
- Continual pre-training of CodeLlama-7B on OpenWebMath with 10B tokens refined by ProX matches Llemma-7B, also pre-trained from CodeLlama-7B, trained on 200B tokens.
- We release "ProX" series of pre-training dataset including >600B general pre-training data, and 5B math corpus.

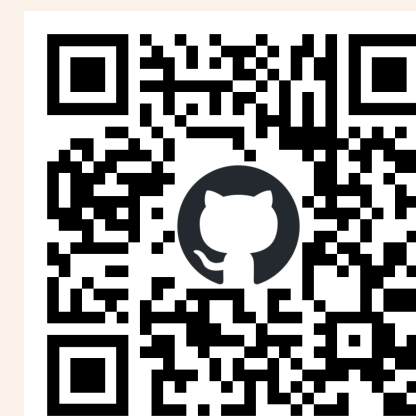## How ProX Works? *Program Design* 👉 *LLM Annotation* 👉 *SFT on Tiny LLMs*

### ProX: Programming Every EXample



**Pre-training Corpora**

[000] Activities | Technical Reports | About W3C | ...
[001] A description of the type of authorized interactions a can have with a resource. Examples include read, write, execute, add, modify, and delete. ClickHere
......
[108] http://t.co/XUQZZRvn3i
[109] Search English terms starting with the letter:

*Noisy & Mixed-quality*

**Base Model**

① *Adapt* 🦙

**Refining Model**

② Generate

**ProX programs**

*Chunk-level Program*
normalize('ClickHere', '')
remove_lines(0, 0)
remove_lines(108, 109)

*Doc-level Program*
keep_doc()

③ 🐍 Execute

**Refined Corpora**

[001] A description of the type of authorized interactions a can have with a resource. Examples include read, write, execute, add, modify, and delete.
......  ✓

*Clean & High-quality*

🥹 Due to visa restrictions, the authors are not able to attend in person. If you have any questions, please also feel free to reach out via the QR codes!
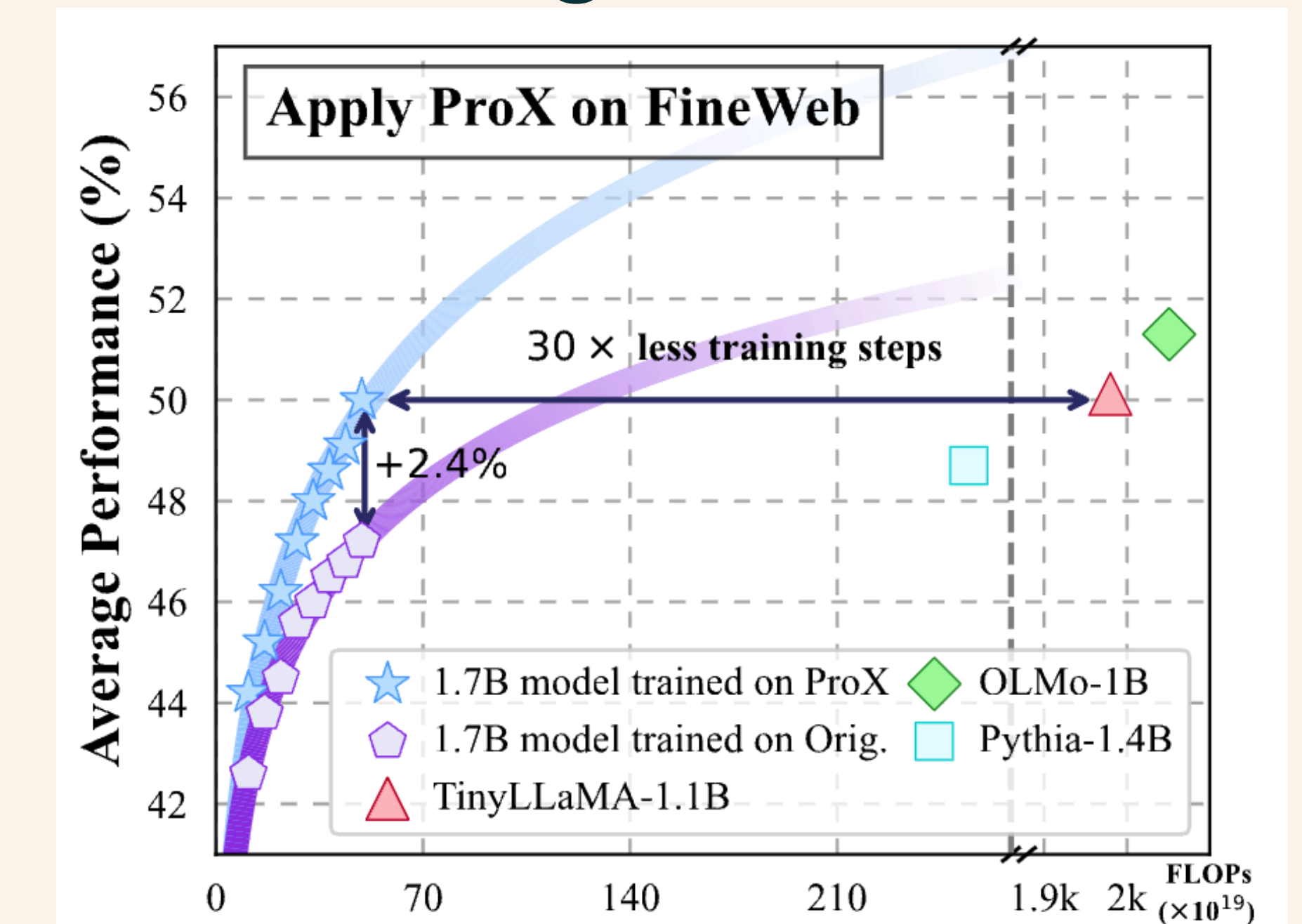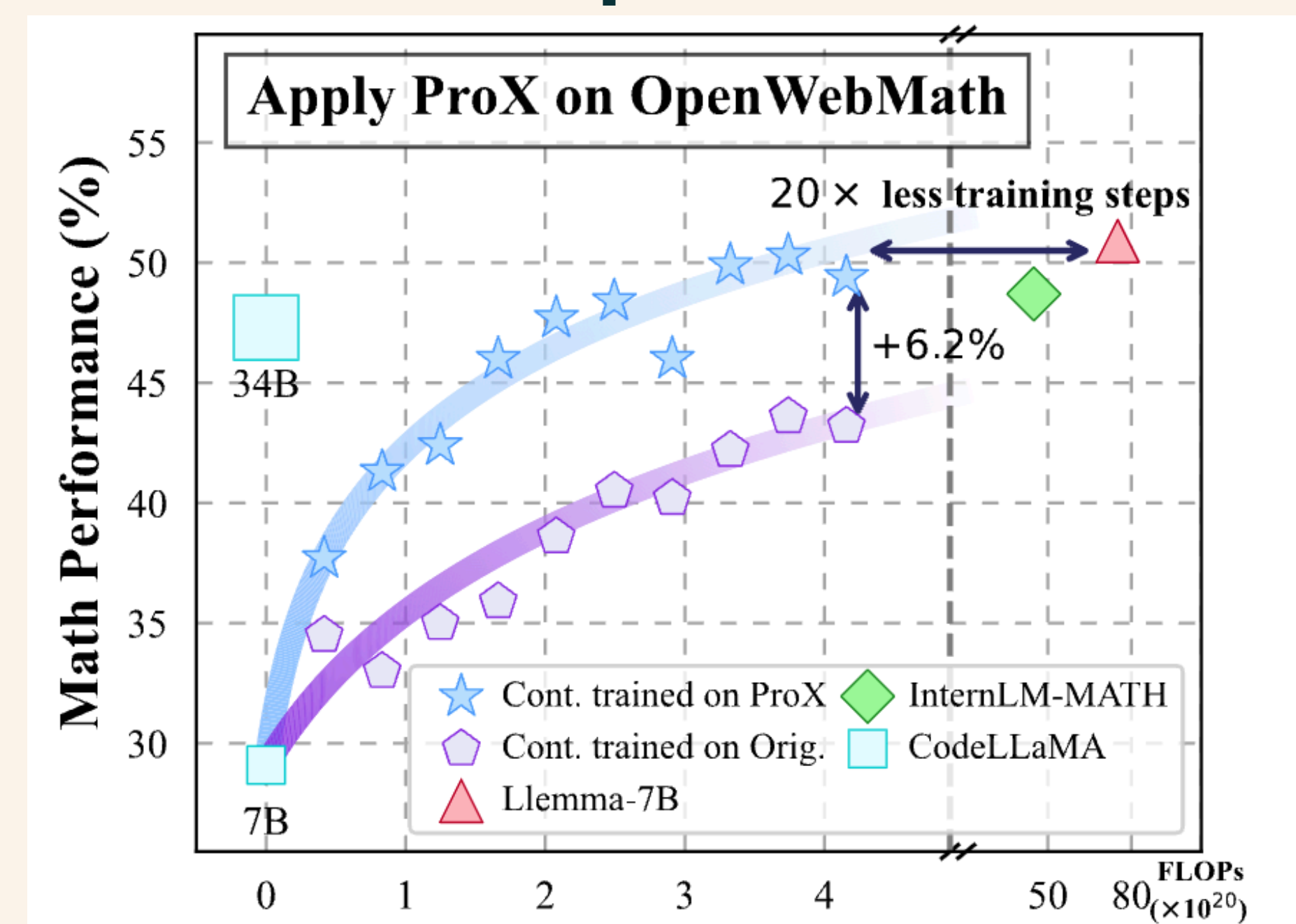
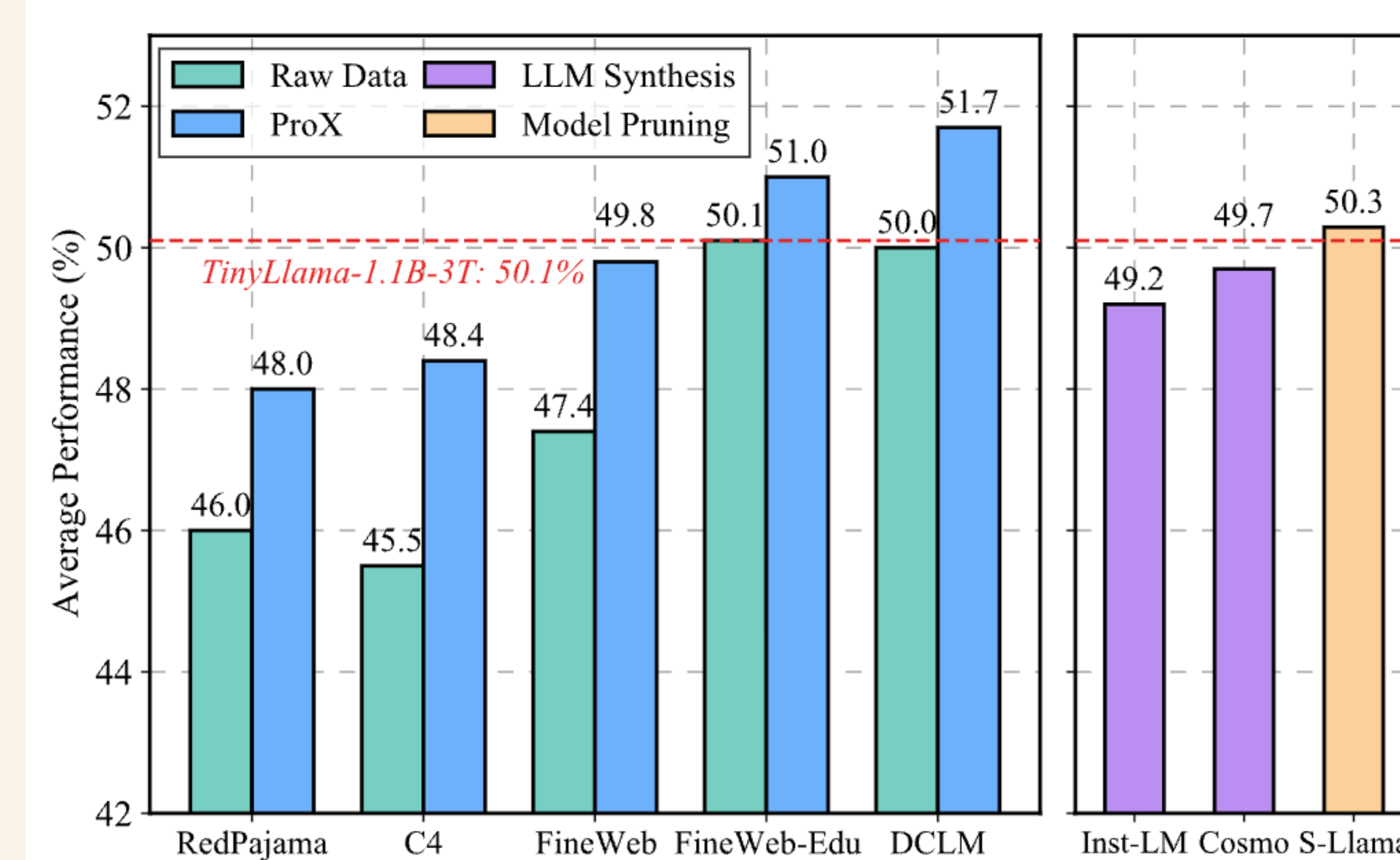Github   Huggingface   X.COM

## Experimental Results

### It works on general domain.



Apply ProX on FineWeb
30 × less training steps
+2.4%
- 1.7B model trained on ProX
- 1.7B model trained on Orig.
- TinyLLaMA-1.1B
- OLMo-1B
- Pythia-1.4B

### It works on specific domain.



Apply ProX on OpenWebMath
20 × less training steps
+6.2%
34B
7B
- Cont. trained on ProX
- Cont. trained on Orig.
- Llemma-7B
- InternLM-MATH
- CodeLLaMA

### And also better than .......



*TinyLlama-1.1B-3T: 50.1%*

Raw Data, ProX, LLM Synthesis, Model Pruning

RedPajama: 46.0, 48.0
C4: 45.5, 48.4
FineWeb: 47.4, 49.8
FineWeb-Edu: 50.1, 51.0
DCLM: 50.0, 51.7
Inst-LM: 49.2
Cosmo: 49.7
S-Llama: 50.3

### With less overall compute!



Train FLOPs, Infer FLOPs

0.3: 0.42, 0.43
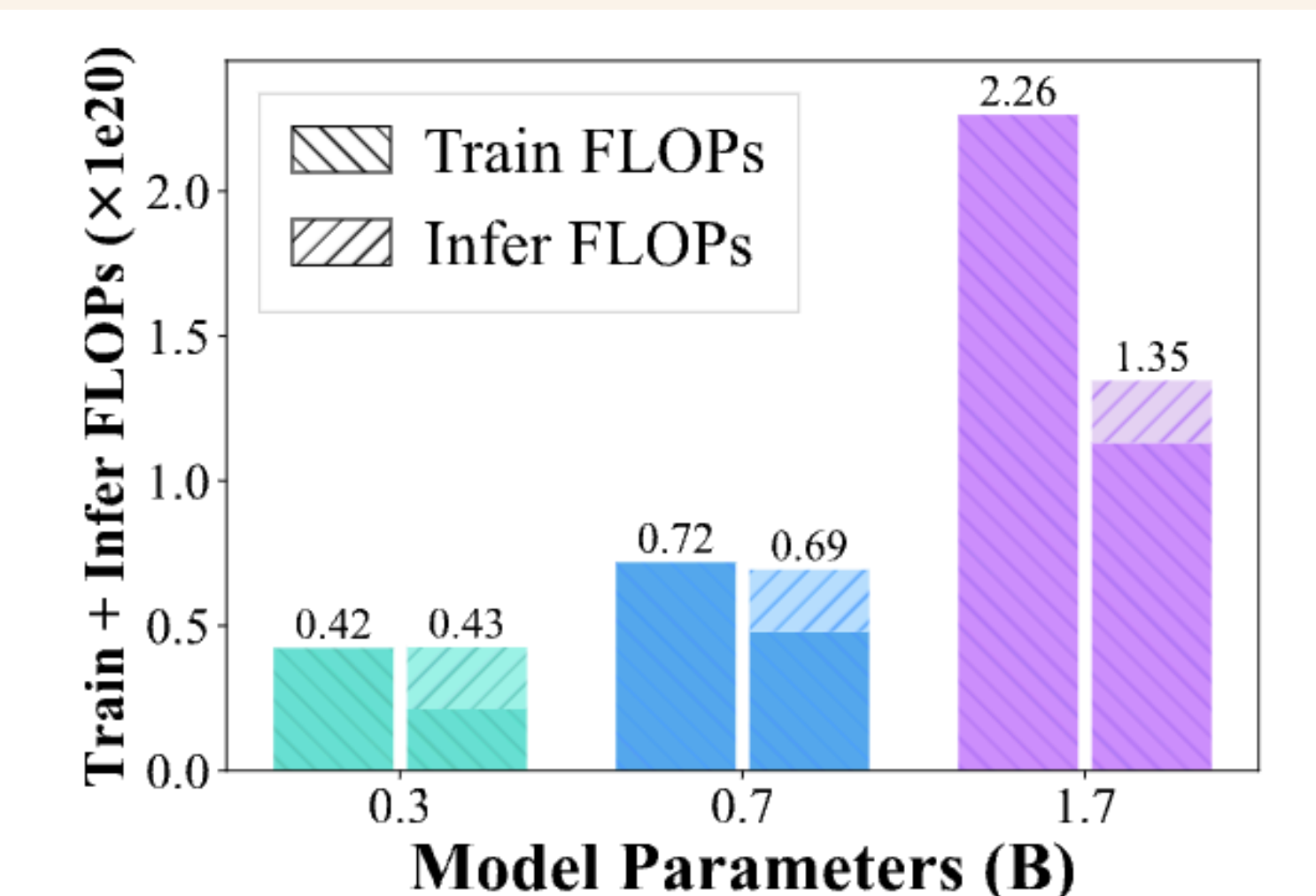0.7: 0.72, 0.69
1.7: 2.26, 1.35

Figure 8: FLOPs comparison for comparable downstream performance with/without PROX refining: 0.3B(Avg.Perf = 40.5), 0.7B (41.6), and 1.7B (42.9).[2]