

# **Scaling Hardness-Impact Training:** INVESTIGATING DATA ORDERING WITH TINYLLAMA<sup>[1]</sup> Hannah Lee<sup>†</sup> Sidharth Lakshmanan<sup>†</sup> Kevin Farhat<sup>†</sup> Chongjiu Gao<sup>†</sup>

#### SUMMARY

- Data ordering by model-based hardness values results in a maximum of 1.38% improvement on finetuned TinyLlama performance on SNLI<sup>[2]</sup>.
- We introduce S-Loss, a technique to scale loss by hardness values during training to

#### RESULTS

Results after finetuning TinyLlama with S-Loss on the SNLI for 2 epochs on 4 GPUs with a per GPU batch size of 8, a learning rate of 4e-4, and AdamW.

	Baseline	CONFIDENCE	NORMALIZED CONFIDENCE	NORMALIZED PERPLEXITY	VARIABILITY	Normalized Variability
TEST ACCURACY	72.36%	<u>74.04%</u>	<u>74.04%</u>	73.56%	73.64%	73.88%

mimic data ordering in a distributed setting, matching/exceeding 1.68% improvement.

### **Research Questions**

- How can we leverage data hardness metrics to improve model performance?
- How can we adapt curriculum learning for a distributed setting?

# **CURRICULUM LEARNING**<sup>[3, 5]</sup>

Ordering finetuning data by hardness showed that ascending difficulty outperforms descending difficulty for all hardness calculations.



![](_page_0_Figure_15.jpeg)

S-Loss

A custom scaling function was created to scale the fine-

	S-Loss I	-unction	10
2.00		Hardness = 0 Hardness = 1	1.0

![](_page_0_Figure_19.jpeg)

#### DATA HARDNESS METRICS

#### **DATASET CARTOGRAPHY**<sup>[4]</sup>

A model-based tool to characterize and diagnose datasets by evaluating various metrics during training. The two prominent metrics are:

- Output Confidence in the true class.
- <u>Variability</u>: The variability of the confidence across epochs.

#### **N-GRAM PERPLEXITY**

tuning Cross-Entropy loss by a coefficient. The coefficient to scale by is defined by the hardness of the sample (normalized) and the current training step (normalized by the total number of steps).

$$\mathscr{L} = \alpha H(p,q) = -\alpha \sum p(x) \log q(x)$$

$$\alpha = \frac{2h_x - 1}{s - 1} + 1 \quad \text{where} \quad h_x \in [0, 1], \quad s \in [0, 1]$$

![](_page_0_Figure_29.jpeg)

Examples	EASY PREMISE	EASY HYPOTHESIS	HARD PREMISE	HARD HYPOTHESIS
Perplexity	Being exposed to hot sun like this can cause skin cancer.	Using sunblock may prevent diseases	A construction worker sweeping	A girl drinks soda
VARIABILITY	A black dog in a grassy field fetching a colourful toy.	A boy threw the frisbee.	A man Rollerblades across a yellow pole at night.	A man is sleeping at his cubicle at work.
Confidence	A boy wields a net in a boat with another in the middle of a lake.	A boy rides in a hot air balloon.	Photographers take pictures of a girl sitting in a street.	The photographer is taking a picture of a boy

An even linear interpolation of unigram, bigram, and trigram perplexity on the concatenated inputs.

## FUTURE WORK

- It is important to see if these results generalize to different datasets and different tasks. For example, with a CausalLM on GSM8k.
- Explore and perform ablation studies on various loss scaling functions and hardness metrics.
- Compare other finetuning optimizations to S-Loss results to gauge performance improvements.

![](_page_0_Figure_36.jpeg)

![](_page_0_Figure_37.jpeg)

#### Indicates equal contribution

[1] TinyLlama <u>https://arxiv.org/abs/2401.02385</u>
[2] SNLI Dataset <u>https://arxiv.org/pdf/1508.05326v1.pdf</u>

[3] Curriculum Learning for Natural Language Understanding <u>https://aclanthology.org/2020.acl-main.542.pdf</u>

[4] Dataset Cartography <u>https://arxiv.org/pdf/2009.10795.pdf</u>
[5] Curriculum Learning for Language Modeling <u>https://arxiv.org/pdf/</u>