



北京大学  
PEKING UNIVERSITY

## Language Models are Super Mario: Absorbing Abilities from Homologous Models as a Free Lunch

---

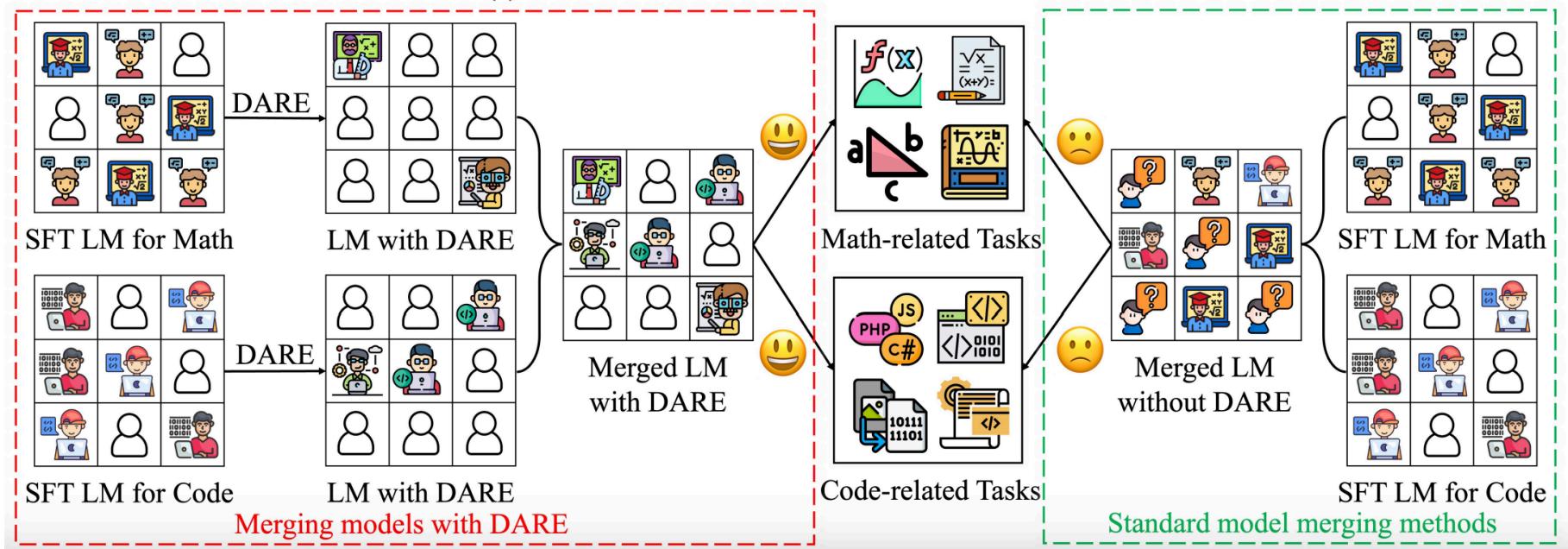
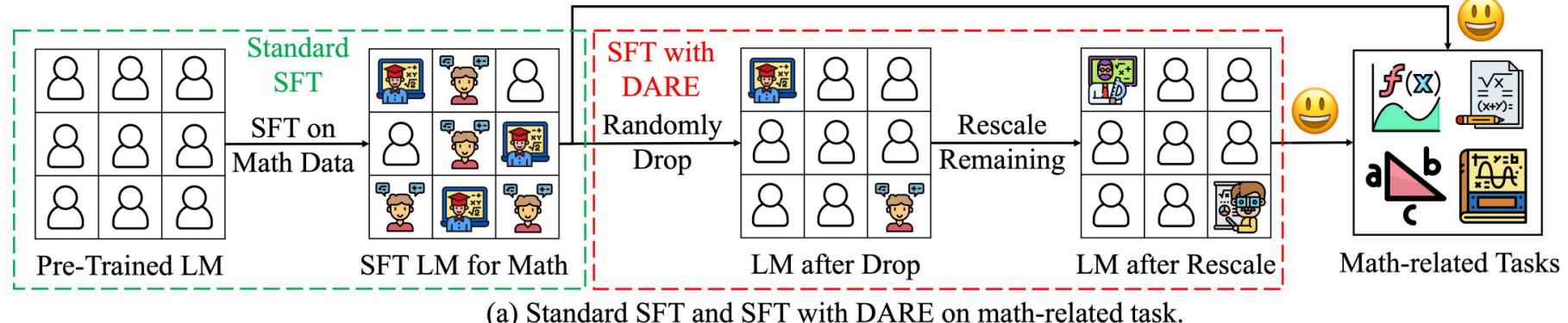


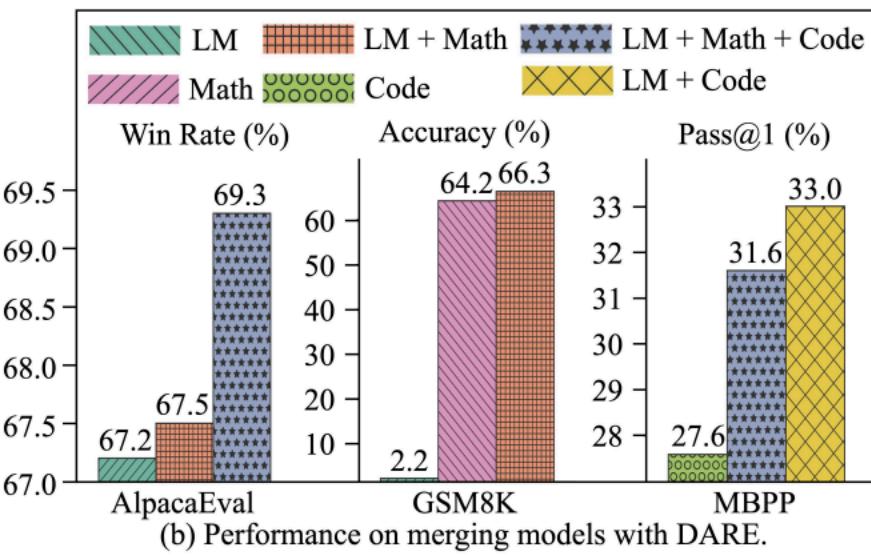
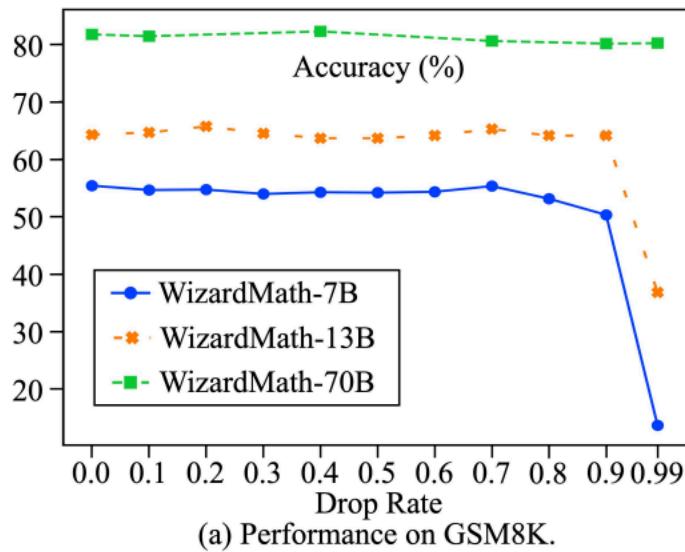
<https://arxiv.org/pdf/2311.03099v2.pdf>

Jiayu Yao

# Framework

Abilities on Math- and Code-related Tasks:





# Lottery Ticket Hypothesis & LoRA

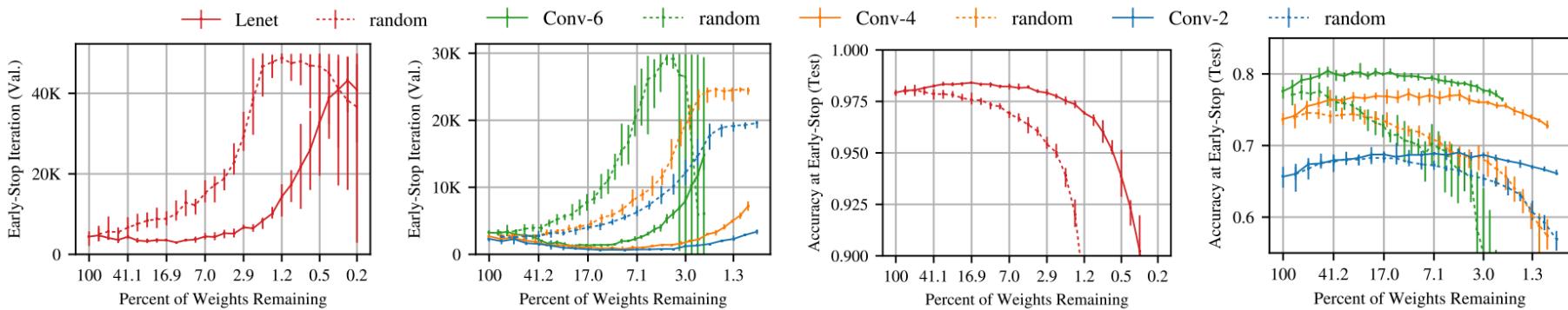


Figure 1: The iteration at which early-stopping would occur (left) and the test accuracy at that iteration (right) of the Lenet architecture for MNIST and the Conv-2, Conv-4, and Conv-6 architectures for CIFAR10 (see Figure 2) when trained starting at various sizes. Dashed lines are randomly sampled sparse networks (average of ten trials). Solid lines are winning tickets (average of five trials).

$$h = W_0x + \Delta Wx = W_0x + BAx$$

$$\delta^t = \theta_{SFT}^t - \theta_{PRE}$$

$$m^t = Bernoulli(p)$$

$$\tilde{\delta}^t = (1 - m^t) \odot \delta^t$$

$$\hat{\delta}^t = \gamma \tilde{\delta}^t$$

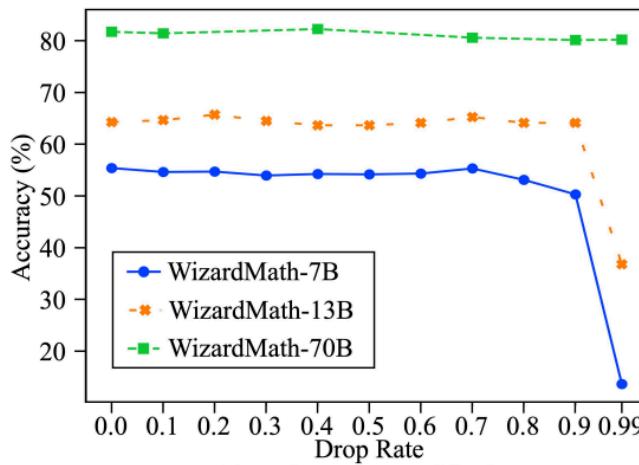
$$\gamma = \frac{1}{1-p}$$

# Proof

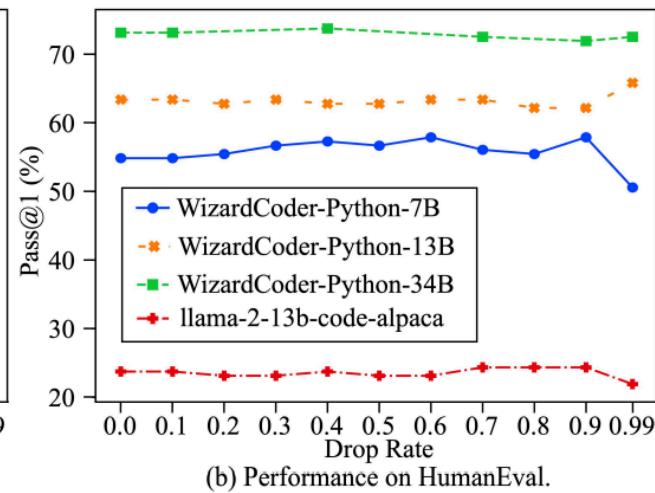
$$\begin{aligned}\mathbb{E}[h_i] &= \mathbb{E}\left[\sum_{j=1}^n (w_{ij} + \Delta w_{ij}) x_j + (b_i + \Delta b_i)\right] \\ &= \sum_{j=1}^n x_j \mathbb{E}[w_{ij}] + \mathbb{E}[b_i] + \sum_{j=1}^n x_j \mathbb{E}[\Delta w_{ij}] + \mathbb{E}[\Delta b_i] \\ &= \sum_{j=1}^n w_{ij} x_j + b_i + \sum_{j=1}^n \Delta w_{ij} x_j + \Delta b_i = h_i^{\text{PRE}} + \Delta h_i\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\hat{h}_i] &= \mathbb{E}\left[\sum_{j=1}^n (w_{ij} + \Delta \hat{w}_{ij}) x_j + (b_i + \Delta \hat{b}_i)\right] \\ &= \sum_{j=1}^n x_j \mathbb{E}[w_{ij}] + \mathbb{E}[b_i] + \sum_{j=1}^n x_j \mathbb{E}[\Delta \hat{w}_{ij}] + \mathbb{E}[\Delta \hat{b}_i] \\ &= \sum_{j=1}^n w_{ij} x_j + b_i + \sum_{j=1}^n x_j ((1-p) \cdot \gamma \cdot \Delta w_{ij} + p \cdot 0) \\ &\quad + ((1-p) \cdot \gamma \cdot \Delta b_i + p \cdot 0) \\ &= h_i^{\text{PRE}} + (1-p) \cdot \gamma \cdot \left(\sum_{j=1}^n \Delta w_{ij} x_j + \Delta b_i\right) \\ &= h_i^{\text{PRE}} + (1-p) \cdot \gamma \cdot \Delta h_i\end{aligned}$$

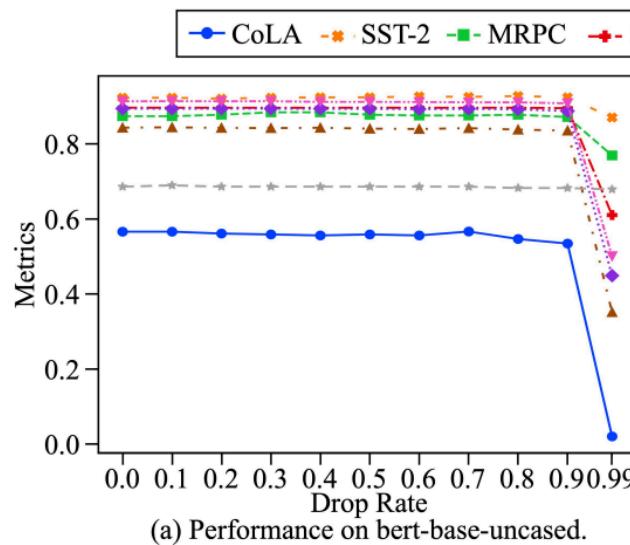
# Experiment



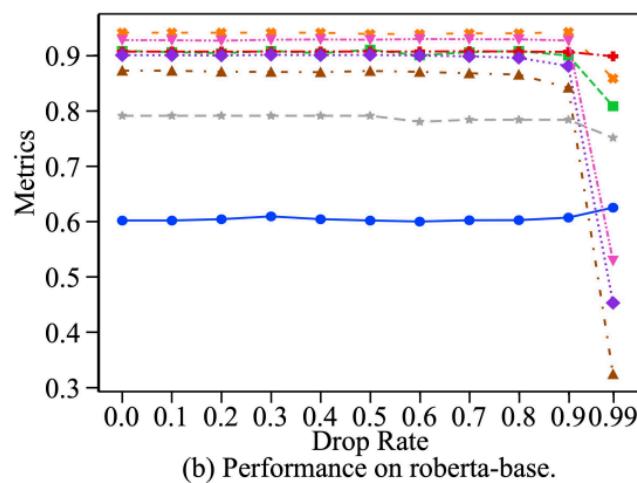
(a) Performance on GSM8K.



(b) Performance on HumanEval.



(a) Performance on bert-base-uncased.



(b) Performance on roberta-base.

# Experiment

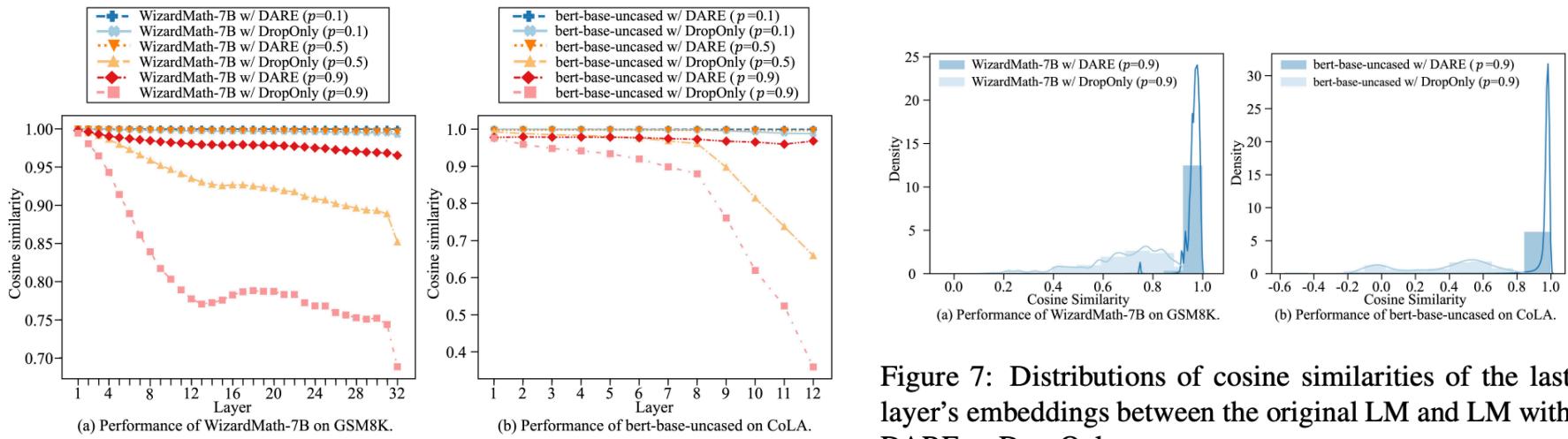
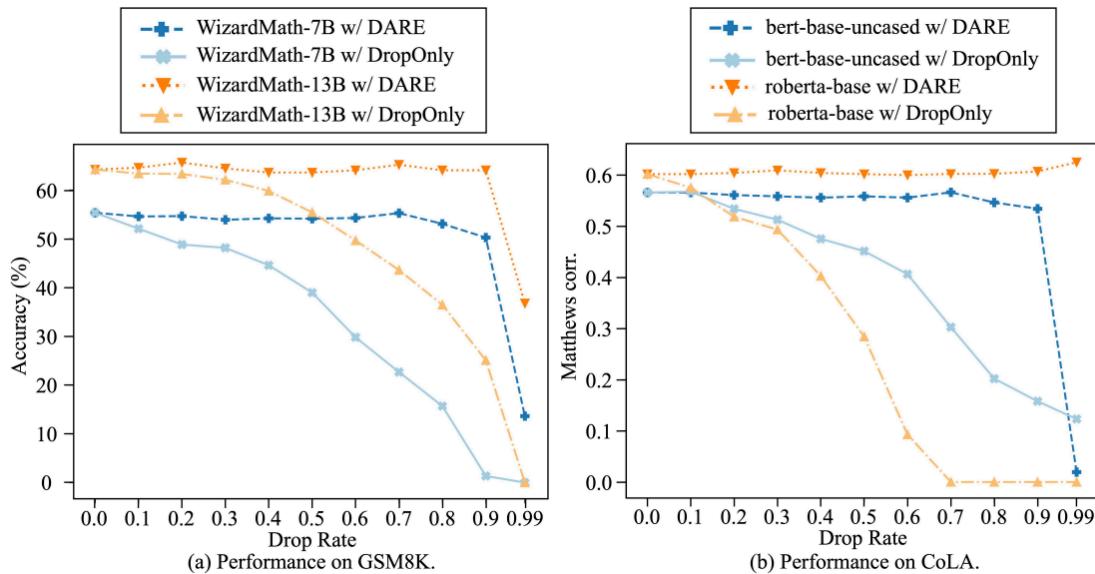
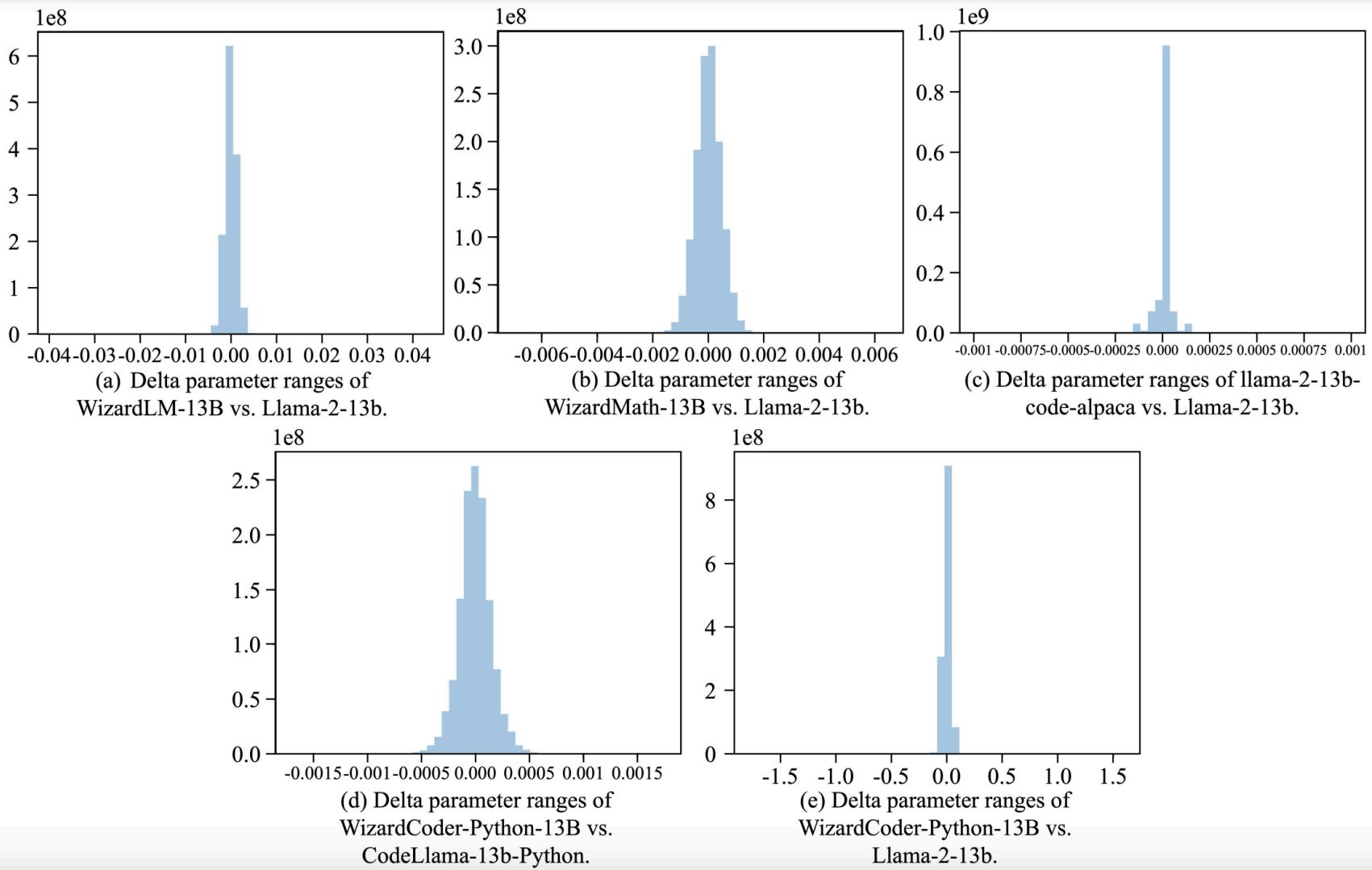


Figure 7: Distributions of cosine similarities of the last layer’s embeddings between the original LM and LM with DARE or DropOnly.



# Experiment



# Experiment

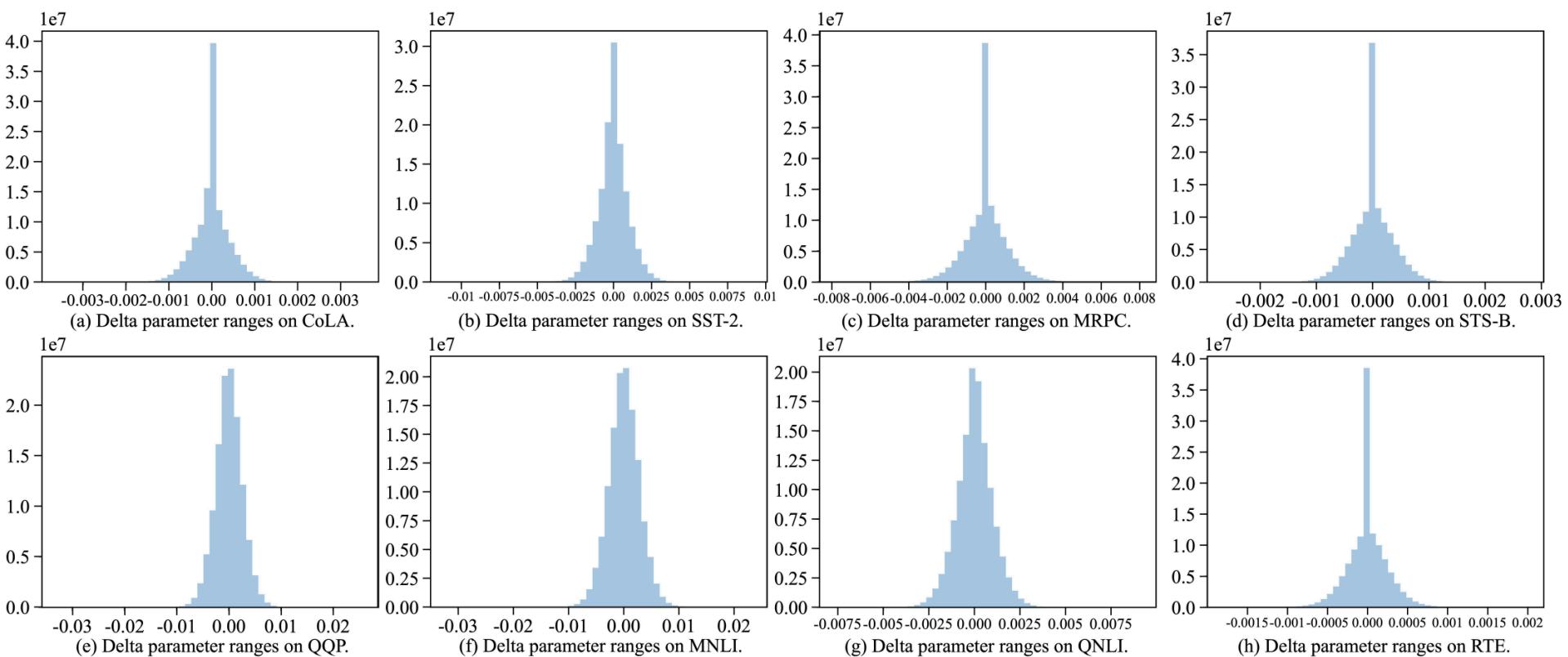


Figure 20: Delta parameter ranges of roberta-base after SFT on GLUE.

# Experiment

Table 1: Performance of merging decoder-based WizardLM-13B (LM), WizardMath-13B (Math), and llama-2-13b-code-alpaca (Code). We use single and mixed colors to denote individual and merged models. The best and second-best results are marked in **bold** and underlined fonts.

Merging Methods	Models	Use DARE	AlpacaEval	GSM8K	MBPP
Single Model	LM	No	67.20	2.20	34.00
	Math	No	/	64.22	/
	Code	No	/	/	27.60
Task Arithmetic	LM	No	67.04	<b>66.34</b>	30.60
	& Math	Yes	<b>67.45</b>	<u>66.26</u>	<u>32.40</u>
	LM	No	<b>68.07</b>	/	<u>32.40</u>
	& Code	Yes	<u>67.83</u>	/	<b>33.00</b>
	Math	No	/	<u>64.67</u>	8.60
	& Code	Yes	/	<b>65.05</b>	<u>9.80</u>
	LM & Math	No	<u>69.03</u>	<u>58.45</u>	<u>29.80</u>
	& Code	Yes	<b>69.28</b>	56.48	<b>31.60</b>

# Experiment

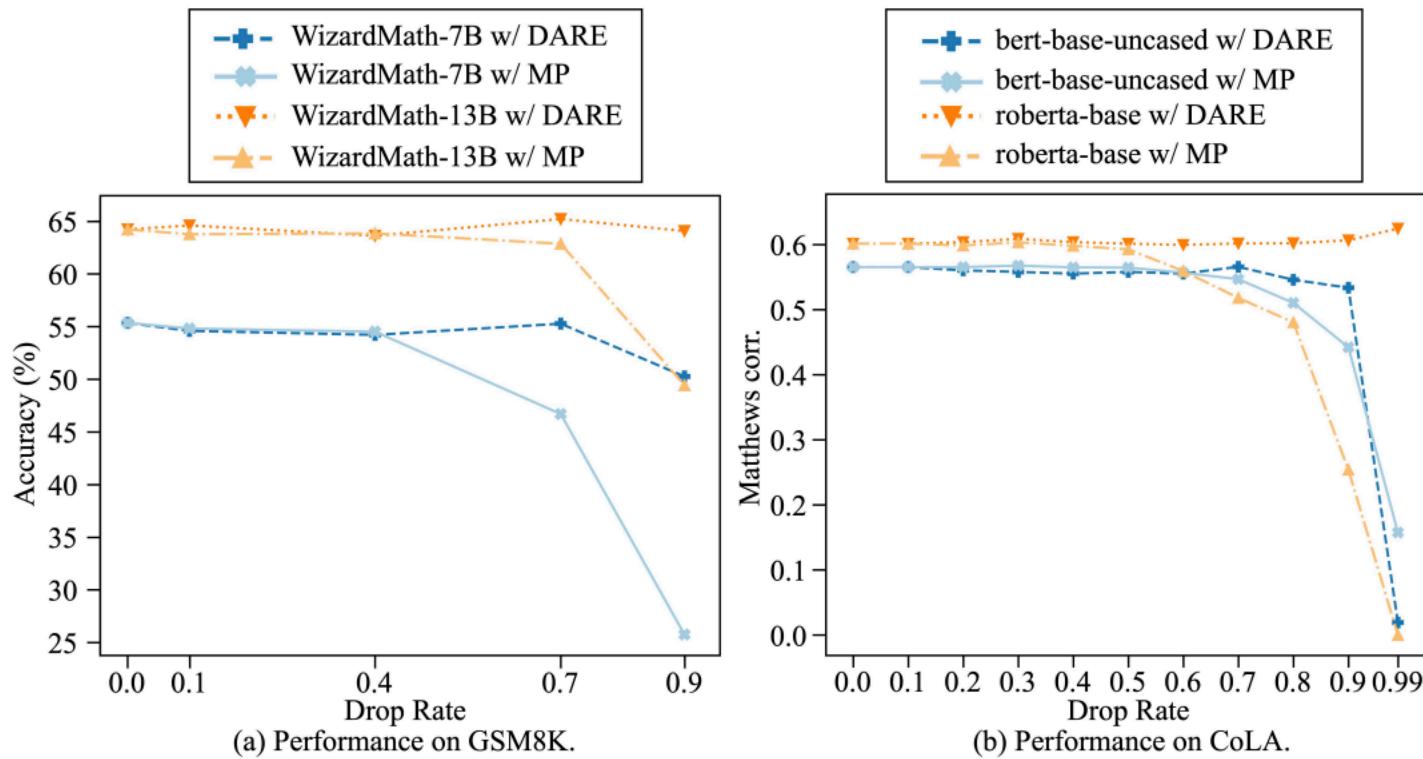


Figure 9: Comparisons between DARE and MP on GSM8K and CoLA on various LMs.

# Discussion

---

- Incremental Learning
- Efficient Tuning
- Subnetwork

# Thanks

---