



# SFT Doesn't Always Hurt General Capabilities: Revisiting Domain-Specific Fine-Tuning in LLMs

Jiacheng Lin<sup>[1]</sup>, Zhongruo Wang<sup>[2]</sup>, Kun Qian<sup>[2]</sup>, Tian Wang<sup>[2]</sup>, Arvind Srinivasan<sup>[2]</sup>, Hansi Zeng<sup>[3]</sup>, Ruochen Jiao<sup>[2]</sup>, Xie Zhou<sup>[2]</sup>, Jiri Gesi<sup>[2]</sup>, Dakuo Wang<sup>[6]</sup>, Yufan Guo<sup>[2]</sup>, Kai Zhong<sup>[2]</sup>, Weiqi Zhang<sup>[2]</sup>, Sujay Sanghavi<sup>[4]</sup>, Changyou Chen<sup>[5]</sup>, Hyokun Yun<sup>[2]</sup>, Lihong Li<sup>[2]</sup>

[1] University of Illinois Urbana-Champaign [2] Amazon [3] University of Massachusetts Amherst [4] University of Texas at Austin [5] University at Buffalo [6] Northeastern University



## Insights from the Analysis

Directly down-weight hard-token losses to curb their potential disproportionate influence to general performance degradation



## Token-Adaptive Loss Reweighting (TALR) for Domain Specific SFT

$$\min_{w \in \Delta_n} \sum_{i=1}^n w_i \cdot \ell_i(\theta) + \tau \sum_{i=1}^n w_i \log w_i$$

- assign smaller weights to harder tokens
- it avoid collapsing all weight onto a small subset of tokens, ensuring broader coverage across the sequence.

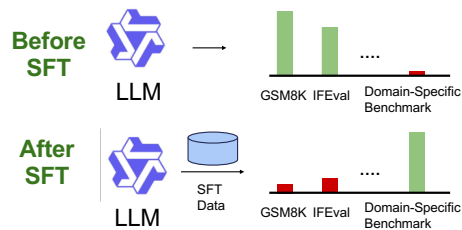
```

Input: Domain dataset  $\mathcal{D}$ , parameters  $\theta$ , learning rate  $\eta$ , temperature  $\tau > 0$ , weight floor  $w_{\min}$ 
Output: Updated parameters  $\theta$ 

foreach training step do
  Sample a mini-batch  $\{(x^{(a)}, y^{(a)}), \dots, (x^{(b)}, y^{(b)})\}_{i=1}^n$  from  $\mathcal{D}$ ;
  Forward pass to obtain token probabilities  $\{p_i\}$  for all supervised tokens in the batch;
  Token NLLs:  $\ell_i \leftarrow -\log p_i$ ;
  Adaptive weights with lower-bound clipping:
     $w_i \leftarrow \exp(-\ell_i/\tau)$ ,  $w_i \leftarrow \max(\text{sg}(w_i), w_{\min})$ 
  Let  $N$  be the number of supervised tokens in the batch;
  Mean (averaged) reweighted loss:
     $\mathcal{L}_{\text{TALR}} = \frac{1}{N} \sum_{i=1}^N w_i (-\log p_i)$ 
  Parameter update:
     $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{TALR}}$ 
end
  
```

Algorithm 1: Token-Adaptive Loss Reweighting (TALR) for Domain-Specific SFT. The  $\text{sg}(\cdot)$  operator denotes stop gradient, meaning that  $w_i$  is treated as a constant during backpropagation to prevent gradients from flowing through the weight computation.

## Recent Papers Claim about Continual Learning: Domain-Specific SFT Harms General Capabilities



## But!! Surprisingly, we found that

Using a smaller learning rate allows domain-specific SFT to achieve a favorable trade-off:

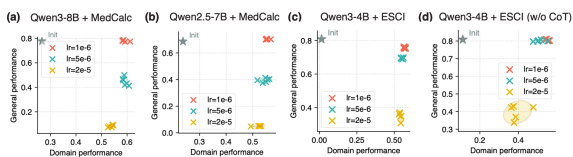
- General-purpose performance degradation is largely mitigated
- Target domain performance is comparable to that with larger learning rates

## Revisiting the Role of Learning Rate

Datasets: two domain-specific datasets, MedCalc and ESCI

## Why choose the two datasets?

- Existing open-source LLMs perform poorly on them
- Thus, domain-specific SFT is motivated to enhance specialized capabilities in domains



## Findings:

- Smaller learning rates achieve a more favorable trade-off.
- Label-only supervision loosens learning rate constraints for Pareto-optimal trade-offs.

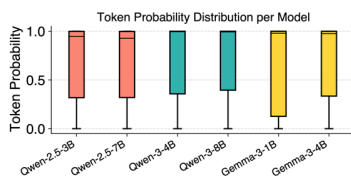
## Why do milder updates preserve general abilities while still enabling strong domain gains?

### Token Analysis

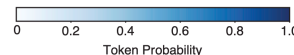
$$p(x_t | x_{\text{prompt}}, x_{<t})$$

We compute

- $x_t$ : each target token during SFT
- $x_{\text{prompt}}$ : prompt tokens
- $x_{<t}$ : previous target tokens



Most tokens in SFT training data pose low learning difficulty.



- The majority of tokens in the target sequence are confidently predicted by the model
- The upper quartiles are tightly clustered near 1.0, and the medians are consistently high

## Theoretical Analysis

We provide a theoretical analysis from the perspective of information theory

**Theorem 3.1.** (Informal) Under certain assumptions, consider fine-tuning on a domain-specific dataset  $\mathcal{D}_2$  with a fixed target domain improvement  $\Delta_* > 0$  (i.e.,  $\Delta L(P_2) \leq -\Delta_*$ ). The general-performance degradation on  $\mathcal{D}_1$ , which is already well modeled by the LLM, admits an upper bound

$$\Delta L(P_1) \leq k_1 \Delta_* + k_2 \Delta_*^2 \lambda$$

where  $\lambda$  is the effective per-step size and  $k_1, k_2$  are constants determined by the model and data. Thus, using smaller steps (smaller  $\lambda$ ) leads to strictly tighter guarantees on general-performance preservation.

**Theorem 3.2.** (Informal) Under certain assumptions, fix a tolerance on general-performance degradation on  $\mathcal{D}_1$  (i.e.,  $\Delta L(P_1) \leq \epsilon_{\mathcal{D}_1}$ ). Then the maximal safe per-step size satisfies  $\lambda_{\max} \propto \frac{\epsilon_{\mathcal{D}_1}}{\sqrt{s}}$ , where  $s$  is the expected number of low-probability tokens per example on  $\mathcal{D}_2$ , defined as tokens whose probabilities under the LLM are below a threshold.

## Experimental Results

Table 1: Comparison of domain and general performance on the MedCalc Benchmark under learning rate 1e-6. Both Standard SFT (with smaller learning rate) and TALR are our contributions, and together they achieve the best overall trade-offs compared with the other baselines.

Method	Qwen2.5-3B		Qwen3-4B		Gemma3-4B		Average	
	Domain	General	Domain	General	Domain	General		
Standard (Ours)	0.4947	0.6202	0.5484	0.7837	0.5587	0.6734	0.5339	0.6924
L2-Reg	0.4904	0.6205	0.4692	0.7964	0.5595	0.6750	0.5064	0.6973
LoRA	0.1261	0.5831	0.1945	0.7640	0.2233	0.1281	0.1813	0.4904
Wise-FT	0.1948	0.6285	0.1428	0.7884	0.2573	0.7635	0.1983	0.7268
FLOW	0.3641	0.5974	0.4768	0.7870	0.5673	0.6914	0.4694	0.6920
TALR (Ours)	0.4806	0.6478	0.4889	0.7880	0.5338	0.7150	0.5011	0.7169

Table 2: Comparison of domain and general performance on the MedCalc Benchmark under learning rate 5e-6. At this larger learning rate, TALR achieves the best overall trade-off by substantially improving general performance while maintaining comparable domain performance.

Method	Qwen2.5-3B		Qwen3-4B		Gemma3-4B		Average	
	Domain	General	Domain	General	Domain	General		
Standard	0.5459	0.3337	0.5782	0.5425	0.5507	0.2655	0.5583	0.3805
L2-Reg	0.5406	0.3470	0.5782	0.5591	0.5471	0.2796	0.5553	0.3952
LoRA	0.1734	0.5670	0.2367	0.7571	0.3864	0.1241	0.2655	0.4827
Wise-FT	0.3584	0.5869	0.3815	0.7531	0.4638	0.5929	0.4012	0.6443
FLOW	0.5266	0.4419	0.5819	0.5599	0.5900	0.3476	0.5528	0.4498
TALR (Ours)	0.5066	0.5490	0.5834	0.6138	0.5351	0.3427	0.5417	0.5018

## Training Dynamics

