

Optimization Inspired Few-Shot Adaptation for Large Language Models

Boyan Gao, Xin Wang, Yibo Yang, David A. Clifton

Motivation

- Limitation of existing few-shot adaption methods:
 - In-context learning (ICL) methods:
 - Instable performance to the format of prompts
 - No learnable parameters to capture the complicated features
 - PEFT (LoRA) based methods:
 - Overfitting to few-shot data
- In this work, we ask:
 - *For few-shot adaptation, how can we develop an efficient method that avoids overfitting to few-shot data, as commonly observed in PEFT, while also overcoming ICL's lack of learnable parameters and extra inference cost?*

Method

- Reframing the forward pass as a *learnable preconditioned gradient descent* (PGD) parameterised by **LayerNorm**:

$$Z_{t+1} = Z_t - P_t \nabla \mathcal{L}(Z_t), \quad P_t = \Gamma_t \cdot \frac{1}{\sigma_t}$$

- Optimization with **fast convergence**:

$$\mathcal{J}(P) = \sum_{t=1}^{T-1} \frac{\|Z_t - Z_{t+1}\|}{\|Z_t - Z_{t-1}\|}$$

- Improving **generalization** through sharpness-aware regularization:

$$\mathcal{H}(P) = \text{tr}(P_t \nabla^2 \mathcal{L}(Z_t) P_t^T)$$

Preconditioned GD in LLM

- LLMs learn layer-wise gradient descent in the **forward pass**

$$\begin{aligned} Z_{t+1} &= Z_t - \eta \nabla \mathcal{L}(Z_t) \\ \text{s.t. } f_t(Z_t) &= -\eta \nabla \mathcal{L}(Z_t) = \text{Attn}(Z_t) \end{aligned}$$

- Reframe this as as preconditioned gradient descent by treating **LayerNorm** as learnable preconditioner $P = \{P_t\}_{t=1}^T$

$$Z_{t+1} = Z_t - \Gamma_t \cdot \frac{\nabla \mathcal{L}(Z_t) - \mu_t}{\sigma_t}, \quad \Gamma_t = \text{diag}(\gamma_t)$$

$$Z_{t+1} = Z_t - P_t \nabla \mathcal{L}(Z_t), \quad P_t = \Gamma_t \cdot \frac{1}{\sigma_t}$$

Fast Convergence

- Fast Convergence objective:

$$\mathcal{J}(P) = \sum_{t=1}^{T-1} \frac{\|Z_t - Z_{t+1}\|}{\|Z_t - Z_{t-1}\|}$$

- Step ratio: optimization efficiency and stability described by:

$$\|Z_{t+1} - Z^*\| \leq \rho_t \|Z_t - Z^*\|, \quad \rho_t < 1$$

- Implicitly optimizing the convergence bound of preconditioned gradient descent

Generalization

- Optimizing the **local sharpness** of loss landscape
 - Hutchinson approximation for the **preconditioning Hessian trace**

$$\begin{aligned}\mathcal{H}(P) = \text{tr}(P_t \nabla^2 \mathcal{L}(Z_t) P_t^T) &\approx \frac{1}{\epsilon} \mathbb{E}_{\nu} \left[\nu^T P_t (\nabla \mathcal{L}(Z_t + \epsilon P_t \nu) - \nabla \mathcal{L}(Z_t)) \right] \\ &\approx \frac{1}{\epsilon} \frac{1}{N} \sum_i \left[\nu_i^T P_t (\nabla \mathcal{L}(Z_t + \epsilon P_t \nu_i) - \nabla \mathcal{L}(Z_t)) \right]\end{aligned}$$

The algorithm for this approximation is given in our paper

Experiments

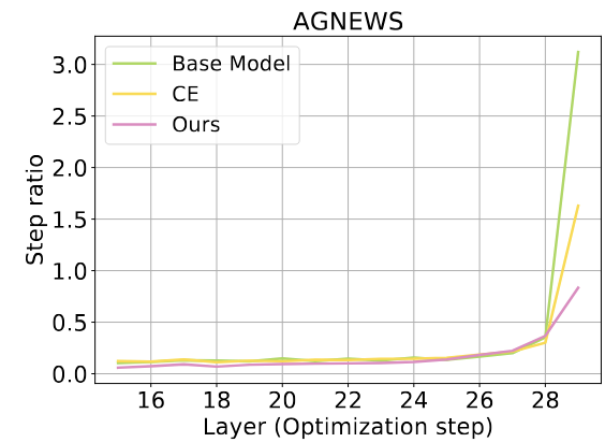
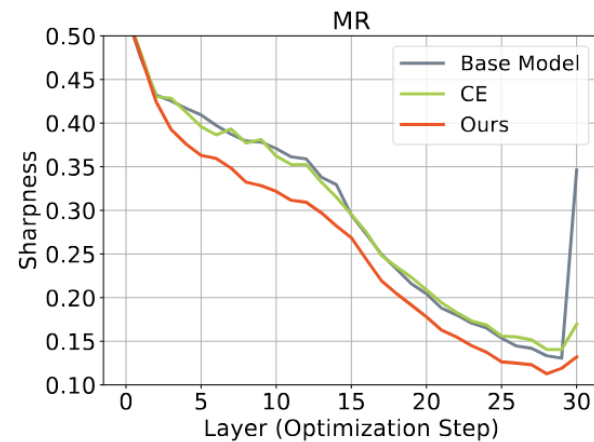
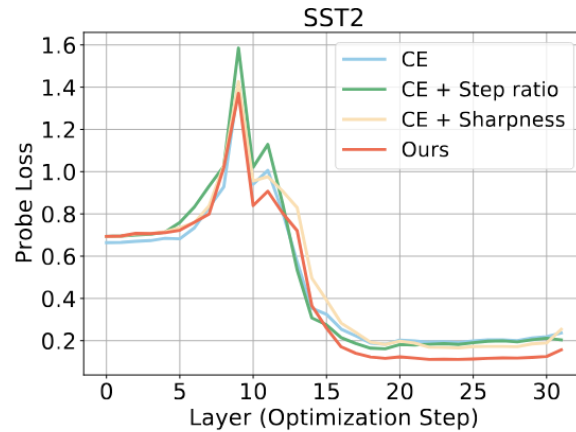
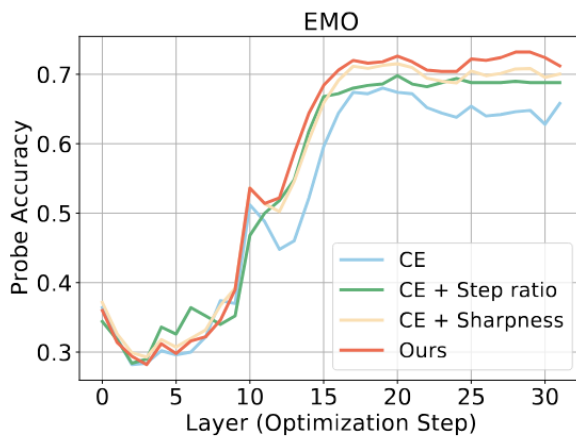
- SOTA on various few-shot benchmarks when competing with the baseline model

Table 1: Comparison between OFA and other baseline algorithms on Llama2-7B and Llama3-8B-Instruct. Mean accuracy and standard deviation across five random seeds are reported. **Best** results are highlighted in bold.

Dataset	SST-2	SST-5	TREC	AGNews	Subj	HateSp18	DBPedia	EmoC	MR
Method	Llama2-7B								
Zero-shot	83.00	27.00	50.00	70.20	51.40	54.20	72.00	41.80	73.60
Few-shot (ICL)	94.44 \pm 1.44	41.72 \pm 3.68	77.32 \pm 4.41	85.68 \pm 2.00	52.56 \pm 3.09	70.24 \pm 5.80	96.64 \pm 0.48	75.48 \pm 1.63	93.24 \pm 0.50
Soft-prompt	56.24 \pm 6.99	24.24 \pm 2.96	55.20 \pm 4.14	78.00 \pm 7.60	57.40 \pm 4.93	59.56 \pm 6.96	74.40 \pm 6.43	35.08 \pm 5.29	54.32 \pm 1.76
Label-anchor	83.32 \pm 5.95	27.68 \pm 4.21	77.48 \pm 3.49	83.72 \pm 1.04	53.00 \pm 2.95	64.52 \pm 8.09	81.40 \pm 3.67	59.12 \pm 10.60	84.40 \pm 5.89
Task-vector	81.44 \pm 4.73	25.96 \pm 0.59	65.68 \pm 1.93	79.68 \pm 4.07	58.56 \pm 4.91	67.68 \pm 3.70	89.48 \pm 2.58	44.64 \pm 3.53	82.32 \pm 5.37
IA3	93.28 \pm 2.29	46.08 \pm 2.11	84.40 \pm 5.99	87.04 \pm 1.97	71.92 \pm 8.08	72.44 \pm 2.59	94.68 \pm 1.09	64.32 \pm 1.95	88.80 \pm 2.28
I2CL	87.68 \pm 2.47	39.12 \pm 2.69	78.56 \pm 5.32	85.48 \pm 1.16	73.84 \pm 3.84	69.88 \pm 5.67	90.16 \pm 1.86	63.72 \pm 1.37	87.68 \pm 2.26
OFA (Ours)	95.84 \pm 0.41	50.36 \pm 3.28	85.92 \pm 1.90	89.00 \pm 1.26	88.40 \pm 4.76	83.04 \pm 3.72	97.72 \pm 0.52	76.60 \pm 2.39	94.36 \pm 1.13
Method	Llama3-8B-Instruct								
Zero-shot	93.00	35.80	71.00	80.40	50.80	67.80	67.40	53.60	86.40
Few-shot (ICL)	96.48 \pm 0.48	46.72 \pm 2.64	79.92 \pm 5.83	89.64 \pm 0.59	57.48 \pm 7.08	52.72 \pm 2.35	97.00 \pm 0.28	65.28 \pm 4.29	93.12 \pm 0.16
Soft-prompt	84.68 \pm 7.71	38.40 \pm 5.68	75.68 \pm 8.17	84.96 \pm 3.80	73.28 \pm 5.41	62.72 \pm 5.54	82.88 \pm 6.45	55.32 \pm 9.74	75.76 \pm 7.71
Label-anchor	93.36 \pm 2.39	40.54 \pm 5.44	78.28 \pm 4.07	84.64 \pm 1.61	54.16 \pm 2.25	69.48 \pm 5.43	87.48 \pm 3.04	59.36 \pm 2.48	88.20 \pm 3.69
Task-vector	94.80 \pm 2.02	56.42 \pm 1.15	79.83 \pm 1.52	89.21 \pm 0.58	76.08 \pm 1.23	67.12 \pm 0.32	79.52 \pm 1.84	57.96 \pm 4.59	86.52 \pm 0.64
IA3	94.32 \pm 0.82	49.24 \pm 2.06	87.60 \pm 3.46	88.36 \pm 1.80	82.04 \pm 7.43	77.20 \pm 4.37	92.56 \pm 1.82	68.04 \pm 2.24	91.76 \pm 0.43
I2CL	90.84 \pm 0.98	48.96 \pm 2.48	79.60 \pm 6.22	88.96 \pm 2.03	81.48 \pm 4.68	65.88 \pm 3.61	91.20 \pm 2.03	64.32 \pm 2.05	88.88 \pm 0.61
OFA (Ours)	97.08 \pm 0.27	58.32 \pm 2.74	89.06 \pm 1.49	91.84 \pm 0.61	92.64 \pm 3.43	89.47 \pm 0.47	97.92 \pm 1.06	79.24 \pm 4.87	94.56 \pm 0.51

Empirical Analysis

- Probe analysis about the properties introduced by our proposed loss function
 - Layer wise-accuracy, loss, sharpness and step ratio



Thanks