



ICLR 2026 Poster



BA-LoRA: Bias-Alleviating Low-Rank Adaptation to Mitigate Catastrophic Inheritance in Large Language Models

Yupeng Chang¹, Yi Chang^{1,2,3}, Yuan Wu^{1,2,*}

¹ School of Artificial Intelligence, Jilin University, ² Key Laboratory of Symbolic Computation and Knowledge Engineering, Jilin University, ³ International Center of Future Science, Jilin University

Background: PEFT can amplify inherited bias

Phenomenon: Low-rank bottlenecks may amplify inherited bias and noise during downstream adaptation.

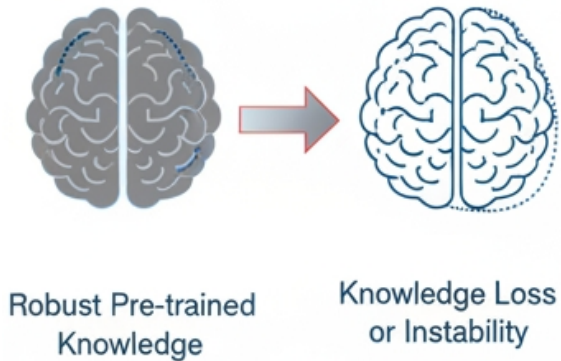


Goal: Mitigate catastrophic inheritance with a lightweight PEFT framework.

Failure Modes of Catastrophic Inheritance

Catastrophic inheritance typically appears in three forms.

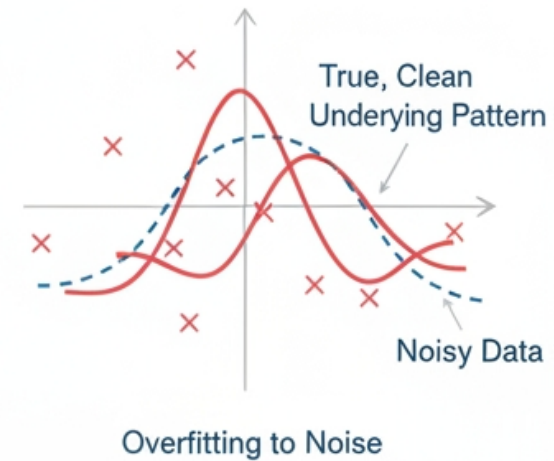
Knowledge Drift



Representation Collapse



Overfitting to Noise



Need: a lightweight PEFT solution that addresses all three failure modes.

BA-LoRA for NLU Tasks

1. Consistency Regularization (CR-NLU)

(CR-NLU): Aligns the fine-tuned model with the pre-trained model via KL divergence.



$$\mathcal{L}_{\text{CR_NLU}} = T^2 \cdot \text{KL}(\text{softmax}(\mathbf{Z}_P/T) \parallel \text{softmax}(\mathbf{Z}_F/T))$$

2. Diversity Regularization (DR-NLU)

(DR-NLU): Reduces inter-class output correlation via covariance regularization.



$$\mathcal{L}_{\text{DR_NLU}} = \frac{1}{D} \sum_{i \neq j} [C(\mathbf{Z}_F)]_{i,j}^2$$
$$C(\mathbf{Z}_F) = \frac{1}{N-1} \mathbf{Z}_{\text{centered}}^T \mathbf{Z}_{\text{centered}}, \quad \text{where } \mathbf{Z}_{\text{centered}} = \mathbf{Z}_F - \bar{\mathbf{Z}}_F$$

3. SVD Regularization (SVDR-NLU)

(SVDR-NLU): Encourages low-rank and robust output representations via SVD.

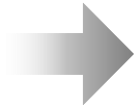


$$\mathcal{L}_{\text{SVDR_NLU}} = -\frac{\sum_{i=1}^k \sigma_i}{\sum_{j=1}^{\min(N,D)} \sigma_j}$$

Overall objective for NLU: $\mathcal{L}_{\text{NLU}} = \mathcal{L}_{\text{task_NLU}} + \lambda_1 \mathcal{L}_{\text{CR_NLU}} + \lambda_2 \mathcal{L}_{\text{DR_NLU}} + \lambda_3 \mathcal{L}_{\text{SVDR_NLU}}$

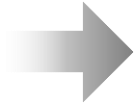
BA-LoRA for NLG Tasks

1. Consistency Regularization (CR-NLG): Aligns the fine-tuned model with the pre-trained model via KL divergence.



$$\mathcal{L}_{\text{CR_NLG}} = T^2 \cdot \frac{1}{M} \sum_{i=1}^M \text{KL}(\mathcal{P}_P(y_i | \mathbf{x}; T) \| \mathcal{P}_F(y_i | \mathbf{x}; T))$$

2. Diversity Regularization (DR-NLG): Promotes diversity within the Top-K candidate tokens.



$$\mathcal{L}_{\text{DR_NLG}} = \frac{1}{M} \sum_{i=1}^M \sum_{v \in \mathcal{V}_{\text{top-K}}^{(i)}} \mathcal{P}'_F(v | \mathbf{h}_i) \log \mathcal{P}'_F(v | \mathbf{h}_i)$$

3. SVD Regularization (SVDR-NLG): Encourages low-rank and robust output structures via randomized SVD.



$$\mathcal{L}_{\text{SVDR_NLG}} = -\frac{\sum_{i=1}^k \tilde{\sigma}_i}{\|\mathbf{Z}_{\text{valid}}\|_F}$$

Overall objective for NLG: $\mathcal{L}_{\text{NLG}} = \mathcal{L}_{\text{task_NLG}} + \lambda_1 \mathcal{L}_{\text{CR_NLG}} + \lambda_2 \mathcal{L}_{\text{DR_NLG}} + \lambda_3 \mathcal{L}_{\text{SVDR_NLG}}$

Main Results

Goal: Compare BA-LoRA against strong PEFT baselines on both NLG and NLU benchmarks.

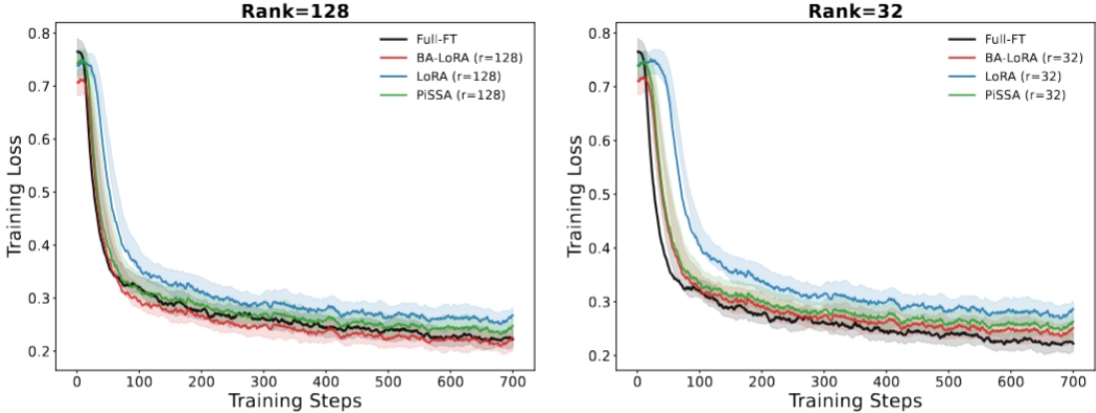
Setup: LLaMA-2-7B for NLG and DeBERTa-v3-base for NLU.

Takeaway: BA-LoRA achieves the best overall performance across both domains. BA-LoRA converges faster and reaches lower training loss.

NLG Benchmarks

Methods	GSM8K	MATH	HumanEval	MBPP	MT-Bench	Avg
Full FT	48.9±0.49	7.48±0.22	20.52±0.29	23.64±0.38	4.85±0.09	21.08
LoRA	42.68±0.54	5.92±0.15	16.80±0.38	21.51±0.43	4.60±0.14	18.30
AdaLoRA	41.95±0.90	6.24±0.38	18.10±0.46	20.19±0.71	4.79±0.18	18.25
DoRA	41.77±0.74	6.20±0.48	16.86±0.54	21.60±0.49	4.48±0.14	18.18
MiLoRA	43.09±1.16	6.31±0.39	17.55±0.24	20.22±0.37	4.50±0.17	18.33
LoRA+	47.84±0.39	7.21±0.49	20.07±0.38	23.69±0.29	5.11±0.06	20.78
LoRA-FA	40.25±0.46	5.66±0.47	15.91±0.41	20.01±0.32	4.67±0.12	17.30
LoRA-GA	50.47±0.98	7.13±0.44	19.44±0.45	23.05±0.40	5.04±0.10	21.03
PiSSA	51.48±0.34	7.60±0.18	19.48±0.45	23.84±0.46	4.92±0.07	21.46
CorDA	53.90±0.56	8.52±0.27	21.03±0.37	24.15±0.44	5.15±0.09	22.55
CorDA++	55.03±0.52	8.95±0.37	21.76±0.39	24.74±0.47	5.64±0.12	23.22
BA-LoRA	55.86±0.35	9.47±0.52	23.58±0.25	36.86±0.31	5.11±0.05	26.18

Optimization Dynamics



NLU Benchmarks

Methods	#Params	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg
Full FT	184M	90.34±0.18	96.33±0.11	89.95±1.07	71.43±0.72	94.24±0.10	92.11±0.28	83.75±1.81	91.04±0.48	88.65
BitFit	0.1M	89.54±0.29	94.68±0.11	87.95±1.33	67.31±0.49	92.45±0.17	88.72±0.45	79.12±0.39	91.63±0.37	86.43
HAdapter	1.22M	90.23±0.07	95.38±0.06	89.97±0.27	68.73±0.27	94.31±0.29	91.99±0.28	84.76±0.39	91.58±0.13	88.37
PAdapter	1.18M	90.42±0.36	95.49±0.10	89.71±0.35	69.04±0.10	94.38±0.26	92.15±0.43	85.53±0.18	91.69±0.13	88.55
LoRA	1.33M	90.71±0.16	94.79±0.16	89.85±0.21	70.05±0.34	93.94±0.09	92.07±0.48	85.43±0.09	91.67±0.29	88.56
DoRA	1.27M	90.48±0.10	95.85±0.08	91.04±0.15	71.03±0.18	94.21±0.37	92.34±0.16	86.19±0.25	91.92±0.38	89.13
AdaLoRA	1.27M	90.87±0.08	96.18±0.43	90.81±0.40	71.64±0.12	94.68±0.46	92.37±0.35	87.78±0.36	91.97±0.43	89.54
PiSSA	1.33M	90.47±0.44	95.81±0.45	91.48±0.49	72.27±0.29	94.41±0.41	92.21±0.26	87.14±0.08	91.93±0.25	89.47
BA-LoRA	1.33M	91.26±0.49	96.25±0.09	92.11±0.55	75.46±0.62	95.35±0.14	93.63±0.52	88.58±0.73	92.71±0.38	90.67

Why Does BA-LoRA Help?

Noisier pre-training:

BA-LoRA shows larger gains when inherited noise is stronger.



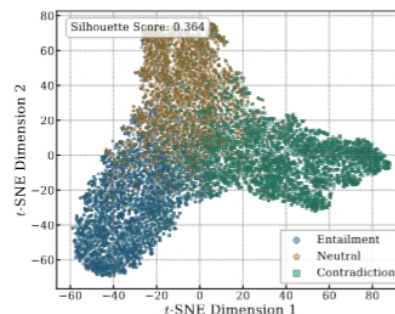
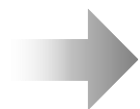
Model	Methods	MNLI	SST-2	CoLA	QNLI	MRPC	Avg
RoBERTa-base	LoRA	85.63 \pm 0.01	94.03 \pm 0.02	62.40 \pm 0.71	91.37 \pm 0.97	87.98 \pm 0.23	84.28
	PiSSA	85.72 \pm 0.40	93.64 \pm 0.13	67.28 \pm 0.59	91.40 \pm 0.54	88.11 \pm 0.24	85.23
	BA-LoRA	86.59\pm0.58	94.83\pm0.45	67.91\pm0.21	92.28\pm0.37	90.07\pm0.32	86.34
T5-base	LoRA	85.30 \pm 0.04	94.04 \pm 0.11	69.35 \pm 0.05	92.96 \pm 0.09	68.38 \pm 0.01	82.01
	PiSSA	85.75 \pm 0.07	94.07 \pm 0.06	74.27 \pm 0.39	93.15 \pm 0.14	76.31 \pm 0.51	84.71
	BA-LoRA	86.91\pm0.48	95.20\pm0.29	80.19\pm1.03	94.12\pm0.32	83.43\pm0.71	87.97

Gain over PiSSA: +1.11

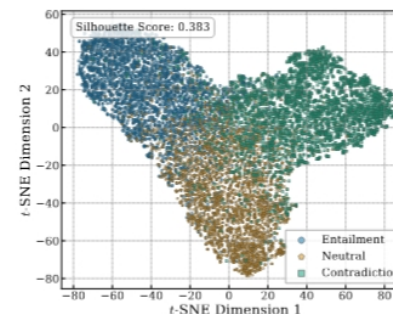
Gain over PiSSA: +3.26

Imbalanced data:

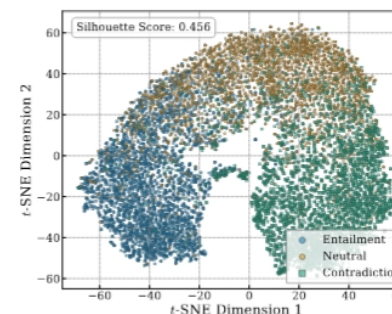
BA-LoRA preserves clearer feature separation under severe imbalance.



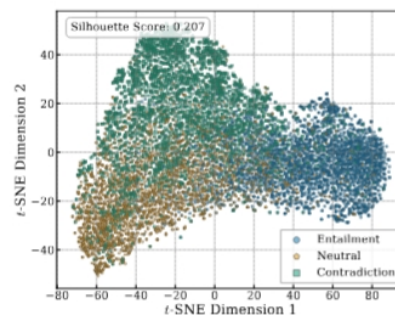
(a) LoRA, Balanced



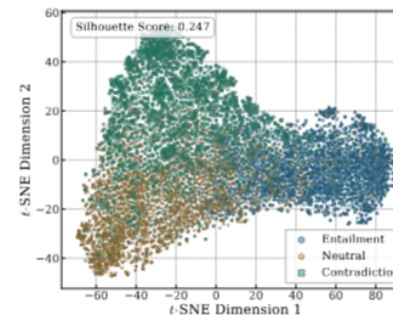
(b) PiSSA, Balanced



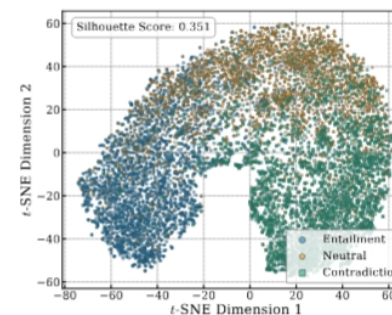
(c) BA-LoRA, Balanced



(d) LoRA, Imbalanced




(e) PiSSA, Imbalanced



(f) BA-LoRA, Imbalanced


Ablation and Efficiency

Ablation Study: All three regularizers contribute positively, and the full BA-LoRA model achieves the best overall performance.



Configuration	GSM8K	MATH	Average of GLUE
Baseline (PiSSA)	51.48 ± 0.34	7.60 ± 0.18	89.47
+ \mathcal{L}_{CR}	54.25 ± 0.59	9.15 ± 0.25	90.18
+ \mathcal{L}_{DR}	53.60 ± 0.46	8.95 ± 0.18	89.85
+ \mathcal{L}_{SVDR}	52.95 ± 0.55	8.70 ± 0.22	89.71
BA-LoRA (Full)	55.86 ± 0.35	9.47 ± 0.52	90.67

Efficiency Analysis: BA-LoRA delivers strong gains with only modest additional memory and training time, achieving a favorable performance-cost trade-off.



Method	Memory Cost	Training Time	GSM8K
Full FT	>96 GB	>24h	48.9 ± 0.49
LoRA	66.32 GB	4h 31min	42.68 ± 0.54
PiSSA	66.59 GB	4h 17min	51.48 ± 0.34
BA-LoRA	77.34 GB	4h 48min	55.86 ± 0.35



Thank you!