SRLoRA: Subspace Recomposition in Low-Rank Adaptation via Importance-Based Fusion and Reinitialization

Haodong Yang¹ Md Zakir Hossain^{1 2 3} Lei Wang^{2 4}

¹Australian National University ²Data61/CSIRO ³Curtin University ⁴Griffith University

MOTIVATION

- LoRA (Low-Rank Adaptation) is a popular Parameter-Efficient Fine-Tuning (PEFT) method by introducing two low-rank matrices A and B while keeping the pre-trained weights frozen.
- LoRA's update $\Delta W = BA$ lies in a subspace of full fine-tuning update. As a consequence, the fine-tuning performance of LoRA is suboptimal.
- We propose Subspace Recomposition in Low-Rank Adaptation via Importance-Based Fusion and Reinitialization:
 - Identify unimportant components in LoRA based on importance scores.
 - Fuse unimportant LoRA weights back into pretrained weights.
 - Reinitialize unimportant components using 'unused' principal components from SVD of pretrained weights.



Before Training:

• **Perform Singular Value Decomposition (SVD)** on the pretrained weights *W*₀, and initialize the LoRA matrices *A* and *B* with the PiSSA method (Meng et al., 2024).

Training Phase: At the end of each training interval,

• Compute the sensitivity-based importance score (Zhang et al., 2023). We define *:*







DATA

group importance score S_k for the *k*-th components:

$$S_k = \frac{1}{m} \sum_{i=1}^m s(B_{ik}) + \frac{1}{n} \sum_{j=1}^n s(A_{kj})$$

- Split the update $\Delta W = BA$ into two subsets based on their importance scores:
 - \circ *B***₁,** *A***₁: low-importance components -> fused back into frozen weights.**
 - B₂, A₂: high-importance components -> retained for continued training.
- **Fuse** B_1 , A_1 into frozen weights.
- **Reinitialize** B_1 , A_1 using the 'unused' principal components from the SVD of W_0 and subtract initialisation of B_1A_1 from frozen weights.
- **Reset importance scores** for the next training interval.



DATASETS

GLUE Benchmark

- A collection of natural language understanding tasks
- Tasks used: SST-2, MRPC, CoLA, QNLI, RTE, STS-B

Task Name	Metric	Task Description
SST-2	Accuracy	Stanford Sentiment Treebank
MRPC	F1	Microsoft Research Paraphrase Corpus
CoLA	Matthews corr.	Corpus of Linguistic Acceptability
QNLI	Accuracy	Question Natural Language Inference
RTE	Accuracy	Recognizing Textual Entailment
STS-B	Spearman corr.	Semantic Textual Similarity Benchmark

RESULTS

Method	Params/Total Params	SST2	MRPC	CoLA	QNLI	RTE	STSB
LoRA PiSSA	1.33M/184M 1.33M/184M	94.84 94.95	90.78 91.50	69.82 71.58	91.89 93.36	$85.56 \\ 84.84$	91.06 90.62
SRLoRA	1.33M/184M	95.75	90.63	71.18	93.68	85.92	90.59

Table 3: Comparison of LoRA, PiSSA and SRLoRA on DeBERTa-v3-base fine-tuned on selected GLUE tasks

Method	CIFAR100	STL10	MNIST
LoRA SBL oB A	90.06 92 51	99.62	98.89

Table 4: Comparison of LoRA and SRLoRA on ViT fine-tuned on CIFAR100, STL10 and MNIST.

Vision Classification Datasets CIFAR-100 / STL-10 / MNIST

Dataset	Type	# Classes	
CIFAR-100	Object classification	100	
STL-10	Semi-supervised classification	10	
MNIST	Handwritten digit classification	10	

Table 2: Vision Classification Dataset.

DISCUSSION

- SRLoRA enhances training efficiency by recomposing the update subspace, allowing the model to capture higher-rank information and achieve faster convergence.
- However, frequent recomposition during the stable phase may lead to underperformance. This suggests a potential direction for future work: dynamically adjusting the frequency of subspace recomposition to further improve performance.



Figure 6: Training loss curves on the RTE task. We fine-tune DeBERTa-v3-base using the same hyperparameters for LoRA, PiSSA, and SRLoRA. The results show that SRLoRA achieves faster initial loss reduction compared to the other two methods.

References

- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. (2022). LoRA: Low-rank adaptation of large language models. In International Conference on Learning Representations.
- Meng, F., Wang, Z., and Zhang, M. (2024). Pissa: Principal singular values and singular vectors adaptation of large language models. arXiv preprint arXiv:2404.02948.
- Zhang, Q., Chen, M., Bukharin, A., He, P., Cheng, Y., Chen, W., and Zhao, T. (2023). Adaptive bud- get allocation for parameter
 - efficient fine-tuning. In The Eleventh International Conference on Learning Representations.