

WINNER-TAKE-ALL COLUMN ROW SAMPLING FOR MEMORY EFFICIENT ADAPTATION OF LANGUAGE MODEL

ZIRUI LIU^{1*}, GUANCHU WANG^{1*}, SHAOCHEN ZHONG¹, ZHAOZHUO XU², DAOCHEN ZHA¹, RUIXIANG TANG¹, ZHIMENG JIANG³, KAIXIONG ZHOU¹, VIPIN CHAUDHARY⁴, SHUAI XU⁴, XIA HU¹

¹RICE UNIVERSITY, ²STEVENS INSTITUTE OF TECHNOLOGY, ³TEXAS A&M UNIVERSITY, ⁴CASE WESTERN RESERVE UNIVERSITY

{zL105, gw22, hz88, DAOCHEN.ZHA, RUIXIANG.TANG, KAIXIONG.ZHOU, XIA.HU}@RICE.EDU; ZXU79@STEVENS.EDU; ZHIMENGJ@TAMU.EDU; {VIPIN, SXX214}@CASE.EDU

MEMORY BOTTLENECK OF FINETUNING LLMs

- **Forward phase and Backward phase of LLM Finetuning:**

$$\begin{aligned} \text{Forward Pass} \quad & \mathbf{Z} = \text{MatMul}(\mathbf{H}, \mathbf{W}), \\ \text{Backward Pass} \quad & \nabla \mathbf{H} = \text{MatMul}(\nabla \mathbf{Z}, \mathbf{W}^\top), \\ & \nabla \mathbf{W} = \text{MatMul}(\mathbf{H}^\top, \nabla \mathbf{Z}), \end{aligned}$$

where $\text{MatMul}(\cdot, \cdot)$ is the General Matrix Multiplication operation, \mathbf{H} and \mathbf{Z} are the activation and output feature maps, respectively. \mathbf{W} is the weight. $\nabla \mathbf{H}$, $\nabla \mathbf{W}$, and $\nabla \mathbf{Z}$ are the gradient of \mathbf{H} , \mathbf{W} , and \mathbf{Z} , respectively. The activations \mathbf{H} are stored in GPU memory during the forward pass, for calculating the weight gradient $\nabla \mathbf{W}$ in the backward pass.

- **Memory Bottleneck of LLM Finetuning:** Although the model parameters contribute to the memory footprint, activations (e.g., storing \mathbf{H}) are the main memory bottleneck during training. As shown in the right-side figure, for T5 models, activations may take roughly 73 ~ 88% of the total memory, depending on the batch size B and sequential length S .

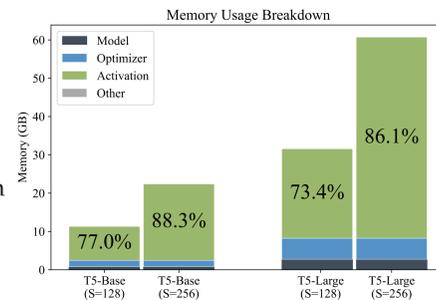


Figure. The GPU memory usage for fine-tuning T5. The batch size is 64 and sequential length is 128 or 256.

WINNER TAKES ALL COLUMN-ROW SAMPLING

- **WTA-CRS Estimator.** WTA-CRS estimator defined in Equation (5) splits the budget k into two parts. Namely, the first part explicitly sums the expectation terms for the largest probability group \mathcal{C} ($|\mathcal{C}| < k$), while stochastically average $k - |\mathcal{C}|$ samples drawn from $\mathcal{D} \setminus \mathcal{C}$ to estimate the remaining terms, up to scale:

$$(\text{WTA-CRS}) \quad \hat{g}(\mathbf{X}, \mathbf{Y}) = \sum_{c \in \mathcal{C}} f(c)p(c) + \frac{1 - \sum_{c \in \mathcal{C}} p_c}{k - |\mathcal{C}|} \sum_{j=1}^{k-|\mathcal{C}|} f(j), \quad i_1, \dots, i_{k-|\mathcal{C}|} \stackrel{\text{i.i.d.}}{\sim} \mathcal{P}^{\mathcal{D} \setminus \mathcal{C}}. \quad (5)$$

- **System Implementation.** As shown in the following right-side figure, a transformer block consists of linear layers, TensorMul, and other operations (e.g., GeLU, Dropout, LayerNorm). TensorMul refers to the multiplication between two four-dimensional tensors. Our WTA-CRS can be applied to the backward pass of Linear-Q, -K, -V, -O, -U, -D, TensorMul-1, and TensorMul-2 (in green), while leaving the forward pass unchanged, as shown in the following left-side figure. The activations of Dropout and GeLU operations (in blue) can be losslessly compressed. The Softmax and LayerNorm operators (in gray) remain unchanged.

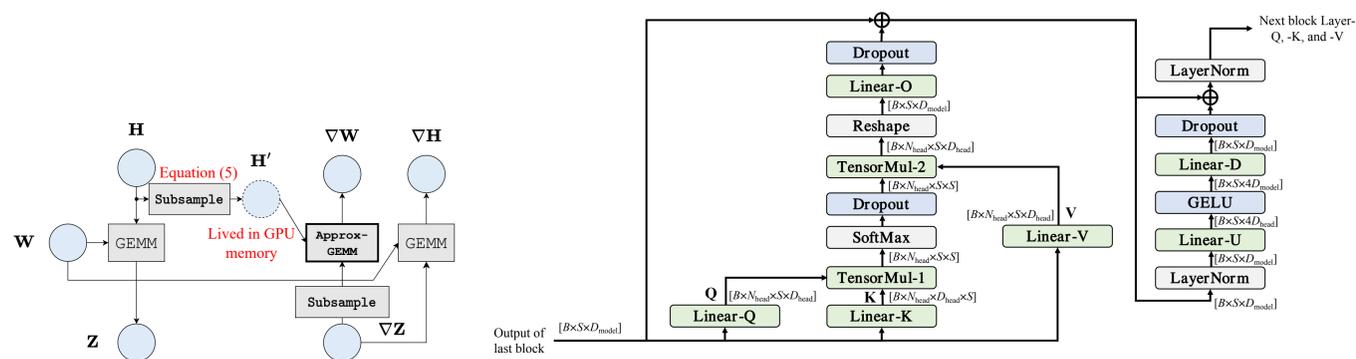


Figure. **Left:** The illustration of how to deploy WTA-CRS to linear layers. **Right:** Application of WTA-CRS to a Transformer block. $B, S, D_{\text{model}}, N_{\text{head}},$ and D_{head} are the batch size, sequence length, hidden size, number of attention heads, and head dimension, respectively. WTA-CRS is applied to the operators in green; the activation maps of operators in blue can be losslessly compressed; and those in gray are not compressed.

MEMORY COST AND ACCURACY ON THE GLUE DATASETS

Table. Peak memory usage (GB) and compression rate of fine-tuning T5-Base and -Large.

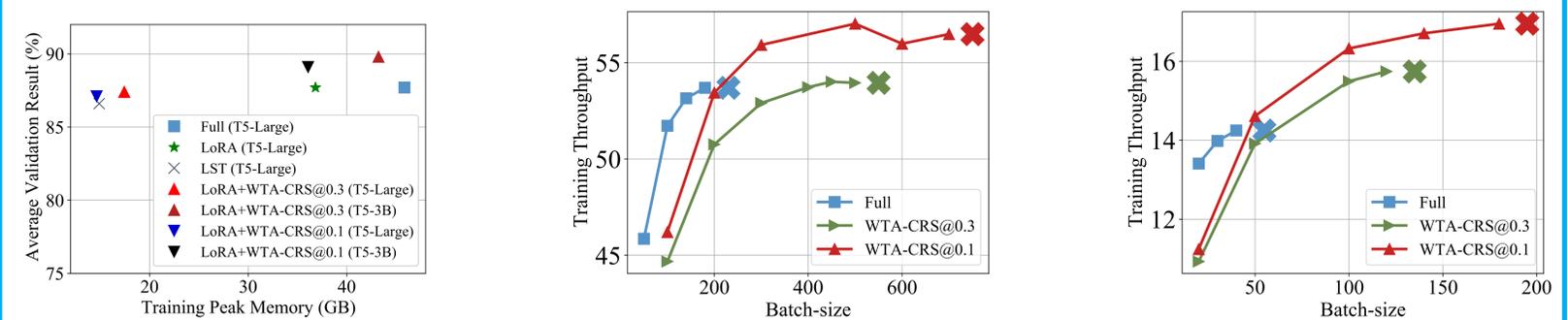
	FP	LoRA	LST	WTA-CRS@0.3	WTA-CRS@0.1	LoRA+WTA-CRS@0.3	LoRA+WTA-CRS@0.1
T5-Base	17.66 (1×)	13.84 (1.3×)	5.50 (3.2×)	8.44 (2.1×)	7.30 (2.4×)	6.50 (2.7×)	5.44 (3.2×)
T5-Large	45.85 (1×)	36.83 (1.2×)	14.85 (3.1×)	21.58 (2.1×)	18.46 (2.5×)	17.44 (2.6×)	14.64 (3.13×)

Table. The GLUE benchmark results with T5 and Bert at different scales.

Model	Method	CoLA	SST-2	MRPC	QQP	MNLI	QNLI	RTE	STS-B	AVG
BERT-Large	Full	66.8±0.31	93.5±0.29	89.5±0.26	88.5±0.03	86.4±0.19	92.1±0.24	72.6±0.36	90.2±0.76	85.0
	LoRA	65.9±0.27	93.8±0.17	90.8±0.37	87.6±0.08	85.9±0.05	92.0±0.2	71.3±0.18	90.3±0.09	84.7
	WTA-CRS@0.3	64.7±0.44	93.5±0.0	89.3±0.39	88.2±0.04	85.2±0.03	91.9±0.12	73.8±0.54	90.4±0.02	84.6
	LoRA+WTA-CRS@0.3	66.0±0.33	93.3±0.29	89.7±1.32	87.6±0.02	86.0±0.07	91.9±0.14	72.4±0.17	89.7±0.04	84.6
T5-Large	Full	61.3±1.01	96.3±0.0	93.4±0.13	89.7±0.01	89.8±0.07	94.2±0.05	85.3±0.17	91.8±0.08	87.7
	LoRA	63.3±0.26	96.4±0.14	93.5±0.16	88.5±0.03	89.5±0.05	94.3±0.07	84.2±0.68	91.7±0.13	87.7
	LST	59.9±0.77	95.8±0.06	91.8±0.08	88.4±0.01	88.7±0.05	94.2±0.02	82.5±0.18	91.4±0.07	86.6
	WTA-CRS@0.3	60.9±1.18	96.3±0.25	93.6±0.57	89.3±0.04	89.5±0.12	94.1±0.03	84.4±0.34	91.3±0.05	87.4
T5-3B	LoRA+WTA-CRS@0.3	62.9±1.19	96.2±0.05	93.6±0.47	88.3±0.02	89.2±0.08	94.0±0.07	83.9±0.95	91.3±0.03	87.4
	LoRA	70.1±0.37	96.8±0.29	94.0±0.27	89.9±0.0	91.0±0.14	95.6±0.05	85.9±0.36	92.9±0.08	89.5
	LoRA+WTA-CRS@0.3	71.4±0.35	96.4±0.06	94.6±0.39	90.0±0.05	91.0±0.06	95.6±0.12	86.3±0.36	92.9±0.09	89.8

- Under similar memory budget, WTA-CRS achieves higher accuracy than other methods, improving down-streaming task performance.
- WTA-CRS achieves a superior trade-off between accuracy and memory usage compared to baselines. Specifically, WTA-CRS has negligible accuracy drop, while the peak memory usage is reduced by $2.1 \times \sim 2.7 \times$ (when combined with LoRA).

EXPERIMENTAL RESULTS



- WTA-CRS achieves better accuracy-memory trade-off than state-of-the-art memory-efficient tuning methods, e.g., LST and LoRA.
- WTA-CRS enables faster training speed under the same hardware. On the T5-Large model, WTA-CRS@0.1 shows $1.08 \times$ higher training throughput; on the T5-3B model, WTA-CRS@0.3 and WTA-CRS@0.1 achieve $1.14 \times$ and $1.21 \times$ higher training throughput, respectively.

ACKNOWLEDGMENTS

The authors thank the anonymous reviewers for their helpful comments. The work is in part supported by NSF grants NSF IIS-1849085, IIS-2224843, and NSF Awards 2117439 and 2112606. This work made use of the High Performance Computing Resource in the Core Facility for Advanced Research Computing at Case Western Reserve University.