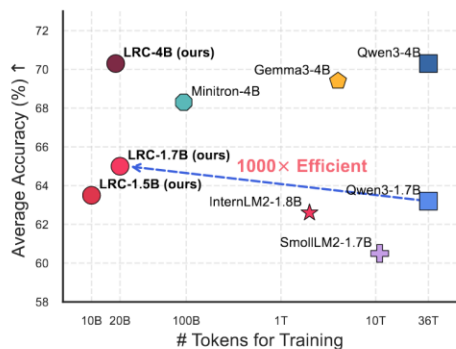# A Token is Worth over 1,000 Tokens: Efficient Knowledge Distillation through Low-Rank Clone

**Spotlight**
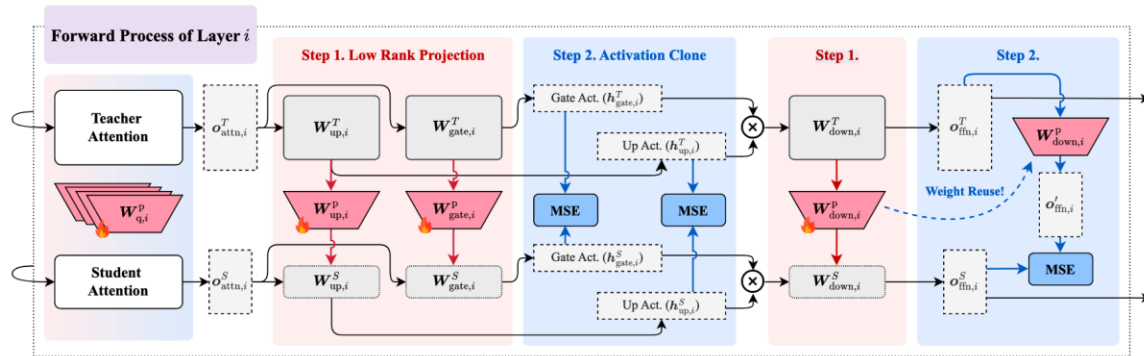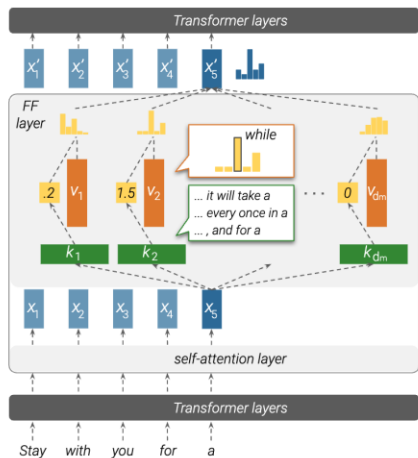
Jitai Hao, Qiang Huang†, Hao Liu, Xinyan Xiao, Zhaochun Ren, Jun Yu†

*HITSZ, Baidu, Leiden, Pengcheng Lab*

LRC surpasses SOTA models trained on trillions of tokens-- while using only **20B** tokens **FROM SCRATCH**, achieving over **1,000x** training efficiency.



## Generate the Student using Low-Rank Projection Instead of Training One!

Attention and normalization modules are omitted. LRC involves two main steps: (1) **Low-Rank Projection**: applying low-rank projection matrices to compress the teacher's weights into a lower-dimensional space, which are then assigned to the student. (2) **Activation Clone**, executing standard forward passes in both models to collect intermediate activations, which are aligned using Mean Squared Error (MSE) loss.

### Insights: Clone the Key-Value Knowledge in FFN.



| Model | Gemma3-4B | Minitron-4B | Qwen3-4B | LRC-4B | LRC-2.7B-B | Sheared-Llama-2.7B-B |
|---|---|---|---|---|---|---|
| Teacher | – | Nemotron4-15B | – | Qwen2.5-7B | Llama2-7B | Llama2-7B |
| # Tokens | 4T | 94B | 36T | 18B | 10B | 50B |
| Dataset | N/A | N/A | N/A | Mixed-2.0 | Redpajama | Redpajama |
| ARC-E | 82.53 | 79.59 | 80.47 | 78.37 | 58.59 | 67.30 |
| ARC-C | 57.08 | 54.35 | 53.58 | 52.47 | 29.61 | 33.58 |
| LogiQA | 33.03 | 30.26 | 33.64 | 34.10 | 29.03 | 28.26 |
| CSQA | 69.37 | 71.09 | 75.76 | 79.28 | 36.36 | 18.92 |
| PIQA | 76.44 | 77.64 | 75.08 | 76.82 | 66.97 | 76.17 |
| WinoG | 69.38 | 65.93 | 65.27 | 67.72 | 62.43 | 65.04 |
| BoolQ | 83.94 | 82.60 | 84.95 | 84.50 | 74.31 | 65.99 |
| SciQ | 95.50 | 96.60 | 95.50 | 95.00 | 85.50 | 91.10 |
| MMLU | 57.58 | 56.77 | 68.38 | 64.41 | 31.20 | 26.56 |
| **Avg. ↑** | 69.43 | 68.31 | 70.29 | **70.30** | **52.67** | 52.55 |

*Main Perf. ↑*



*Ablation ↑*

| Score Type | Teacher | Student |
|---|---|---|
| Original Score | 0.85 | 0.48 |
| Important Neurons Masked | 0.62 (-27%) | 0.33 (-31%) |
| Random Neurons Masked | 0.85 | 0.49 |

*Student FFN clones Teacher's ↑*

| Method | # Tokens/Sec |
|---|---|
| LRC | 84K |
| Sheared Llama (Prune) | 30K |
| Ordinary Training | 146K |
| TinyBERT | 65K |

*Training Efficiency ←*

### Algorithm 1: Overall Procedure of LRC

**Input:** Input token sequence $\mathcal{T}$; number of layers $l$; RMSNorm constant $\epsilon$; teacher's weights $\{W_{m,i}^{\mathrm{T}}\}, W_{\mathrm{emb}}^{\mathrm{T}}, W_{\mathrm{lm}}^{\mathrm{T}}$; low-rank projection matrices $\{W_{m,i}^{\mathrm{P}}\}, W_{\mathrm{emb}}^{\mathrm{P}}, W_{\mathrm{lm}}^{\mathrm{P}}$;

**Output:** Clone loss $\mathcal{L}_{\mathrm{clone}}$;

▷ Step 1: Low-Rank Projection

1 **for** $i = 1$ **to** $l$ **do**
2   **foreach** $m \in \{q, k, v, o, up, gate, down\}$ **do**
3     $W_{m,i}^{\mathrm{S}} \leftarrow W_{m,i}^{\mathrm{T}} W_{m,i}^{\mathrm{P}}$;     ▷ Generate student weights
4 $W_{\mathrm{emb}}^{\mathrm{S}} \leftarrow W_{\mathrm{emb}}^{\mathrm{T}} W_{\mathrm{emb}}^{\mathrm{P}}$; $W_{\mathrm{lm}}^{\mathrm{S}} \leftarrow W_{\mathrm{lm}}^{\mathrm{T}} W_{\mathrm{lm}}^{\mathrm{P}}$;

▷ Step 2: Activation Clone

5 $\mathcal{L}_{\mathrm{clone}} \leftarrow 0$;
6 $h^{\mathrm{T}}, o_{\mathrm{attn}}^{\mathrm{T}}, o_{\mathrm{ffn}}^{\mathrm{T}} \leftarrow \mathrm{Forward}(\mathcal{T}, l, \epsilon, \{W_{m,i}^{\mathrm{T}}\}, W_{\mathrm{emb}}^{\mathrm{T}}, W_{\mathrm{lm}}^{\mathrm{T}})$;   ▷ Get teacher act. dict.
7 $h^{\mathrm{S}}, o_{\mathrm{attn}}^{\mathrm{S}}, o_{\mathrm{ffn}}^{\mathrm{S}} \leftarrow \mathrm{Forward}(\mathcal{T}, l, \epsilon, \{W_{m,i}^{\mathrm{S}}\}, W_{\mathrm{emb}}^{\mathrm{S}}, W_{\mathrm{lm}}^{\mathrm{S}})$;   ▷ Get student act. dict.
8 **for** $i = 1$ **to** $l$ **do**
9   **foreach** $m \in \{q, k, v, gate, up\}$ **do**   ▷ Compute clone loss of interm. states
10     $\mathcal{L}_{\mathrm{clone}} \leftarrow \mathcal{L}_{\mathrm{clone}} + \mathcal{E}(h_{m,i}^{\mathrm{S}}, h_{m,i}^{\mathrm{T}})$;
11   $\mathcal{L}_{\mathrm{clone}} \leftarrow \mathcal{L}_{\mathrm{clone}} + \mathcal{E}(o_{\mathrm{attn},i}^{\mathrm{S}}, o_{\mathrm{attn},i}^{\mathrm{T}} W_{o,i}^{\mathrm{P}}) + \mathcal{E}(o_{\mathrm{ffn},i}^{\mathrm{S}}, o_{\mathrm{ffn},i}^{\mathrm{T}} W_{\mathrm{down},i}^{\mathrm{P}})$;
12 **return** $\mathcal{L}_{\mathrm{clone}}$;

| Model | LRC-1.5B | | |
|---|---|---|---|
| Teacher | Llama3-3B | | |
| # Tokens | 20B | 10B | 10B |
| Dataset | Mixed-2.0 | Mixed-1.0 | Mixed-1.1 |
| **Avg. ↑** | 62.12 | 61.35 | 62.48 |

**Paper** **Github**