# D2O: Dynamic Discriminative Operations for Efficient Long-context Inference in Large Language Models

Zhongwei Wan, Xinjian Wu, Yu Zhang, Yi Xin, Chaofan Tao, Zhihong Zhu, Xin Wang, Siqi Luo, Jing Xiong, Longyue Wang, Mi Zhang

THE OHIO STATE UNIVERSITY

ICLR 2025

## Motivation

- Long-context generative inference in LLMs faces **memory bottlenecks** due to extensive KV cache demands.
- Existing eviction-based methods lead to **context loss and hallucinations**.
- Different layers exhibit varying attention densities; uniform KV cache allocation is **suboptimal**.
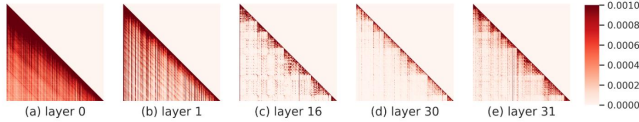
Figure 1: **Attention map density** comparisons of shallow layers (layer 0, 1) and deep layers (layer 16, 30, 31) of LLaMA-2-7B on the GSM8K dataset. We use the mean value of heads for each layer.

## Our Approach: Dynamic Discriminative Operations (D2O)

### Layer-Level Operation:

$$S_l = \alpha_l \cdot S,$$

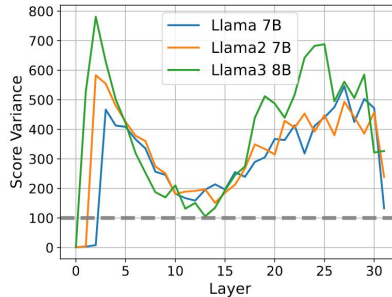$$\text{where } \alpha_l = \frac{\exp(-F_v^l)}{\sum_{l=1}^{L} \exp(-F_v^l)} \cdot L \cdot \rho,$$

Figure 2: **Variances of attention score** across different layers for various models.

- Dynamically allocate KV cache budget ratio using **inverse variance softmax**, prioritizing layers with denser attention maps.
- Figure 2 shows **variances of attention score** across different layers for various models.
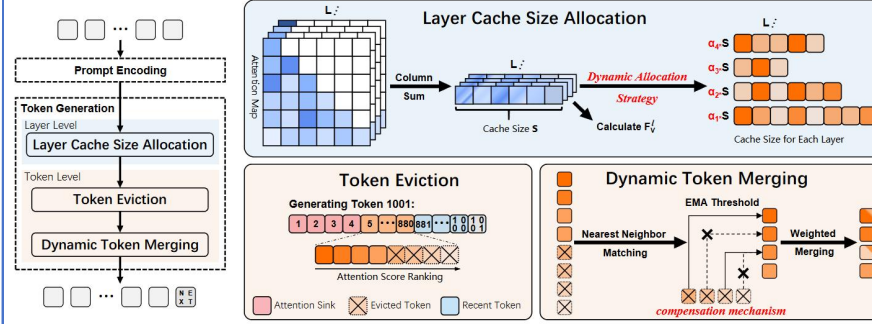
Figure 3: Overview of D2O

## Token-Level Operation:

- **A.** Evict tokens based on cumulative attention with attention sink preservation.
- **B.** Dynamically merge evicted tokens using an **EMA-based similarity threshold** to maintain context relevance.
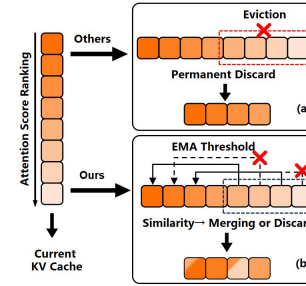
Figure 4: Illustration of dynamic token merging mechanism

**A.**
$$\text{AttnScore} = \begin{cases} \sum_{i=0}^{L_{\text{prompt}}} \mathbf{A}_p[i,:], & \text{if token } i <= L_{\text{prompt}}, \\ \text{Softmax}\left(\mathbf{q}_i \mathbf{K}^\top/\sqrt{D}\right) + \sum_{i=1}^{L_{\text{prompt}}} \mathbf{A}_p[i,:], & \text{otherwise, token generation} \end{cases}$$

$$\mathbf{K}_c = [\mathbf{K}[:T,:], \mathbf{K}[I,:], \mathbf{K}[-M:,:]], \quad \mathbf{V}_c = [\mathbf{V}[:T,:], \mathbf{V}[I,:], \mathbf{V}[-M:,:]],$$
$$\text{and } I = \text{Top}_N\left(\text{AttnScore}[T:-M], N\right),$$

**B.**
$$\mathbf{k}_*^{\text{nearest}} = \underset{j \in I^c}{\text{Argmax}}\left(u_{i,j}\right), \quad \text{where } u_{i,j} = \frac{\mathbf{k}_i^\top \mathbf{k}_j}{\|\mathbf{k}_i\| \|\mathbf{k}_j\|}$$

$$\tau_t = \begin{cases} \frac{1}{L^e}\sum_{i=0}^{L^e} \text{Max}(\mathbf{U}_t[i,:]), & \text{if } t = 0 \text{ for prompt encoding} <= L_{\text{prompt}}, \mathbf{U}_t \in \mathbb{R}^{L^e \times L^c} \\ \beta \text{Max}(\mathbf{U}_t[:]) + (1-\beta)\tau_{t-1} & \text{otherwise}, t > 0 \text{ for token generation}, \mathbf{U}_t \in \mathbb{R}^{L^c} \end{cases}$$

$$\mathbf{k}_{cj} = \mathbf{w}_{cj}\mathbf{k}_{cj} + \sum_{\mathbf{k}_{ei} \in \mathbf{K}_e} \mathbf{w}_{ei}\mathbf{k}_e, \quad \mathbf{v}_{cj} = \mathbf{w}_{cj}\mathbf{v}_{cj} + \sum_{\mathbf{v}_{ei} \in \mathbf{V}_e} \mathbf{w}_{ei}\mathbf{v}_e,$$

$$\mathbf{w}_{cj} = \frac{e}{\sum_{\mathbf{k}_{ei} \in \mathbf{K}_e}\exp(\mathbf{u}_{ij})\mathbf{m}_{ij} + e}, \quad \mathbf{w}_{ei} = \frac{\sum_{\mathbf{k}_{ei} \in \mathbf{K}_e}\exp(\mathbf{u}_{ij})\mathbf{m}_{ij}}{\sum_{\mathbf{k}_{ei} \in \mathbf{K}_e}\exp(\mathbf{u}_{ij})\mathbf{m}_{ij} + e},$$

## Results & Evaluation

Figure 4: Reasoning dataset

Table 1: LongBench benchmarks

Table 2: Throughput Comparison

| Methods | L=50k | L=100k | L=50k | L=100k |
|---|---|---|---|---|
| Full Model | 97.88 | 94.46 | 97.88 | 94.46 |
| | | 4096 | | 8192 |
| StreamingLLM | 58.64 | 47.93 | 62.84 | 51.34 |
| H2O | 79.84 | 69.81 | 82.32 | 72.34 |
| SnapKV | 83.55 | 76.22 | 86.63 | 80.42 |
| CaM | 82.66 | 78.22 | 87.59 | 78.88 |
| D2O | **91.27** | **87.74** | **94.48** | **91.88** |

Table 3: Needle-in-a-haystack

Figure 5: Long sequence modeling

Figure 6: MT bench