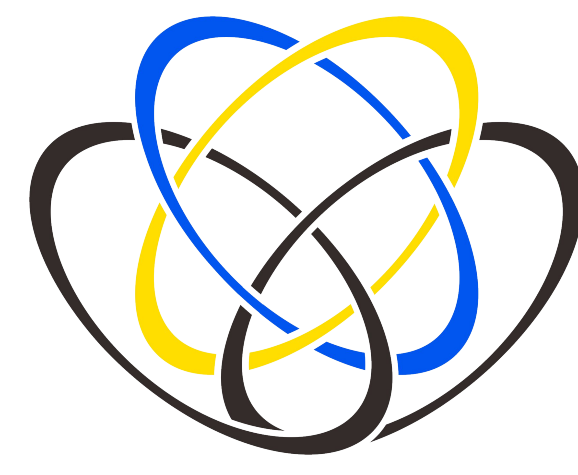
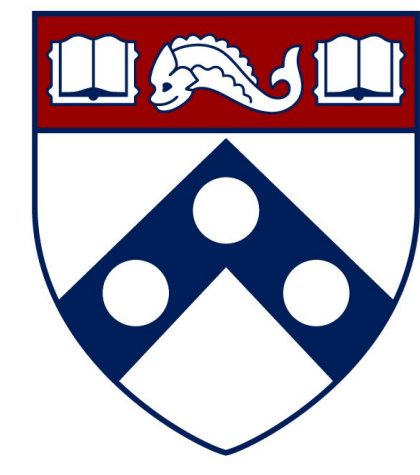


OBCache: Optimal Brain KV Cache Pruning for Efficient Long-Context LLM Inference



Yuzhe Gu¹, Xiyu Liang², Jiaojiao Zhao³, Enmao Diao⁴

¹University of Pennsylvania

²University of Electronic Science and Technology of China

³Duke Kunshan University

⁴DreamSoul



ICML
International Conference
On Machine Learning



Paper



Code

Motivation

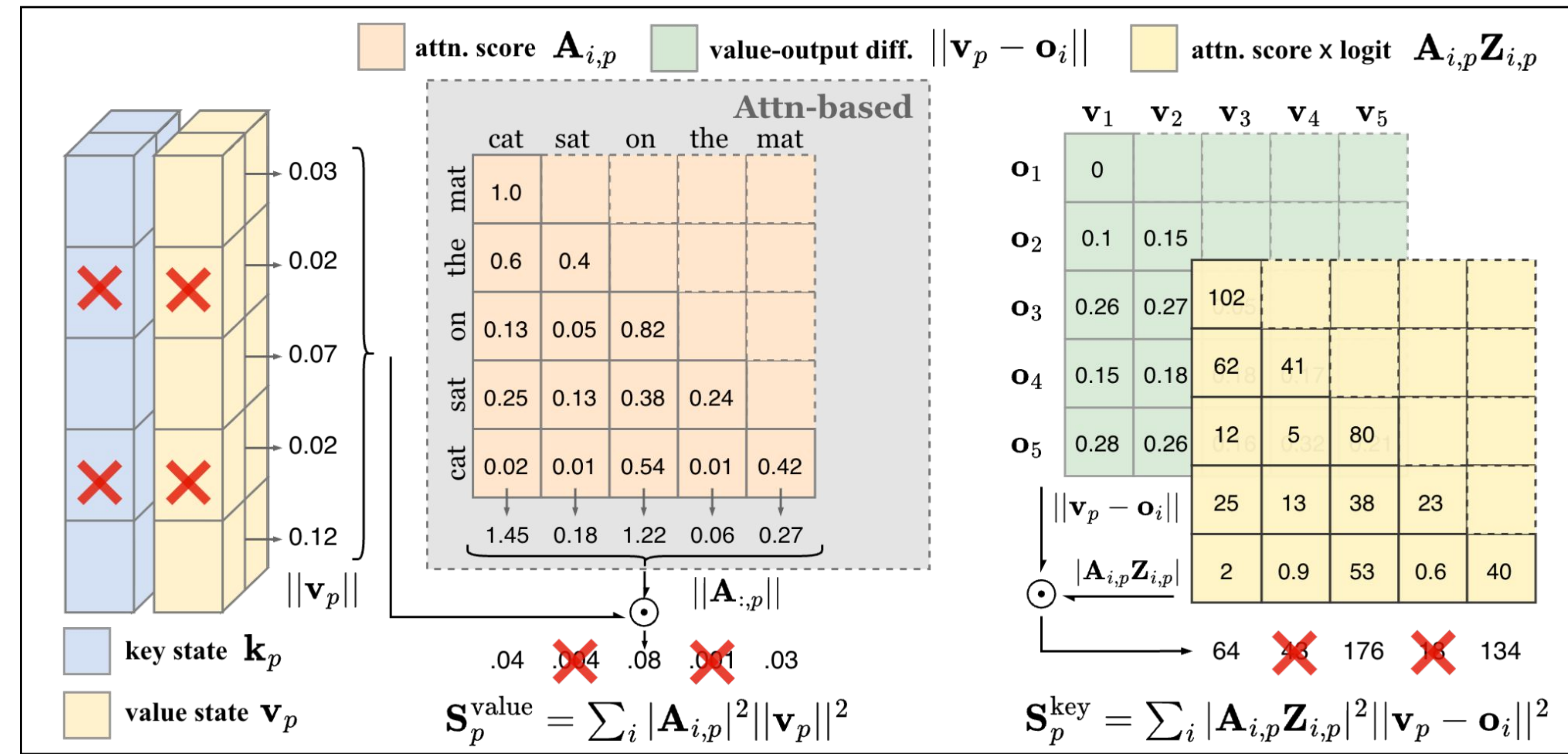
- **Efficiency Bottleneck:** KV cache memory in LLMs grows linearly with sequence length & batch size, creating a key bottleneck in long-context deployment.
- **Attention-Only Heuristics:** Existing KV cache eviction methods mainly estimate token saliency using cumulative attention weights, overlooking value states and the true eviction impact on model outputs.
- **Weak Theoretical Grounding:** Recent KV scoring mechanisms are often empirically driven, lacking a principled theoretical framework for characterizing how token eviction affects model outputs.

Contributions

- **Structured Pruning View:** Building on Optimal Brain Damage (OBD) theory, we formalize KV cache eviction as a layer-wise structured pruning problem, providing a new theoretical framework for cache eviction.
- **Output-Aware Saliency:** We measure token saliency from historical attention output perturbation, deriving three second-order closed-form saliency scores that generalize existing attention-only scores.
- **Plug-and-play Score:** OBCache can be seamlessly integrated into any score-based KV cache eviction pipeline to improve their long-context accuracy, while incurring negligible additional overhead.

OBCache is a Theoretically Grounded Scoring Framework for KV Cache Compression.

Optimal Brain Cache (OBCache)



Token Saliency via Pruning-Induced Eviction Error

$$S_p := \mathcal{L}_{e_p^T[\hat{V} \hat{K}] = 0}(\hat{V}, \hat{K}) = \left\| \text{softmax}\left(\frac{Q_{w:s} \hat{K}^T}{\sqrt{d}}\right) \hat{V} \Big|_{e_p^T[\hat{V} \hat{K}] = 0} - \text{softmax}\left(\frac{Q_{w:s} K^T}{\sqrt{d}}\right) V \right\|_F^2$$

Second-order Taylor Appx. $S_p \stackrel{\text{second order}}{\approx} \frac{1}{2} v_p^T H_{pp}^{vv} v_p + \frac{1}{2} k_p^T H_{pp}^{kk} k_p + v_p^T H_{pp}^{vk} k_p$

OBCache (output-aware) scores

① $S_p^{\text{value}} = \sum_i |A_{i,p}|^2 \|v_p\|^2$

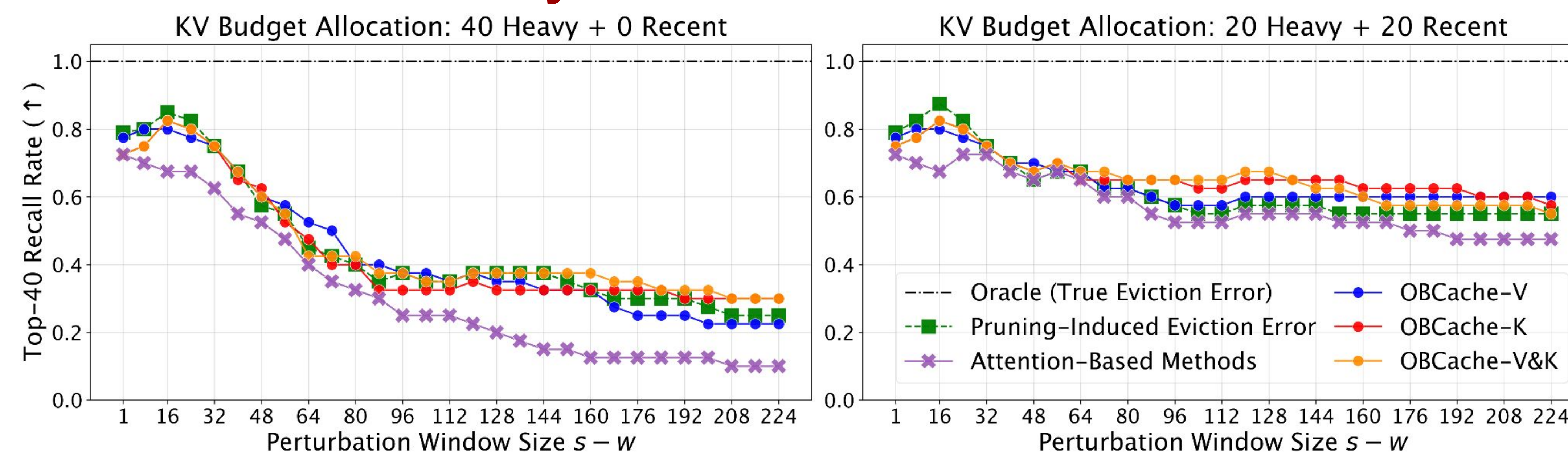
② $S_p^{\text{key}} = \sum_i |A_{i,p} Z_{i,p}|^2 \|v_p - o_i\|^2$

③ $S_p^{\text{joint}} = S_p^{\text{value}} + S_p^{\text{key}} + 2 \sum_i |A_{i,p}|^2 Z_{i,p} (\|v_p\|^2 - v_p^T o_i)$

Attn-only scores

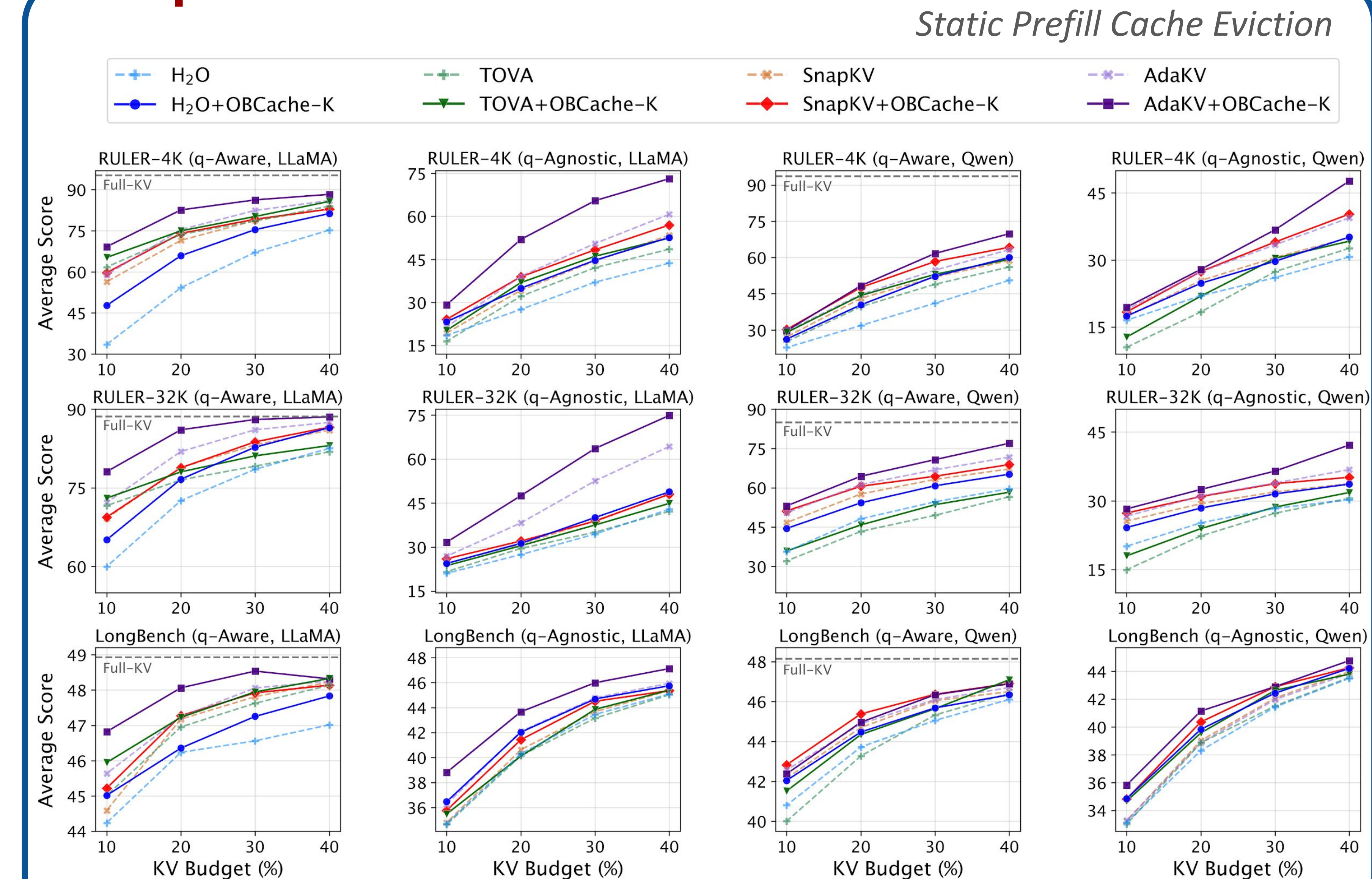
$S_p^{\text{attn}} = \sum_i |A_{i,p}|$

Effectiveness Analysis



- True Eviction Error (unobservable) \rightarrow Pruning-Induced Eviction Error (Proxy Objective)
- Exact Proxy (computationally infeasible) \rightarrow OBCache scores \checkmark Attention-only scores \times

Experimental Results



KV Budget (%)	RULER-4K (Q-Aware)					RULER-4K (Q-Agnostic)					RULER-32K (Q-Aware)					RULER-32K (Q-Agnostic)				
	10	20	30	40	Avg.	10	20	30	40	Avg.	10	20	30	40	Avg.	10	20	30	40	Avg.
Full KV	95.3																			
H ₂ O	33.5	54.2	67.1	75.3	57.5	18.5	27.6	37.0	43.7	31.7	60.0	72.5	78.5	82.5	73.4	21.1	27.5	34.4	42.9	31.5
+ OBCACHE-V	37.7	57.1	68.6	76.0	59.9	19.6	28.2	36.5	44.9	32.3	62.7	73.5	80.4	84.3	75.2	20.0	27.7	36.7	46.0	32.6
+ OBCACHE-K	47.7	65.9	75.4	81.3	67.6	23.3	34.9	44.7	52.7	38.9	65.1	76.6	82.7	86.5	77.7	24.4	31.3	40.1	48.8	36.2
+ OBCACHE-V&K	46.3	67.5	75.6	82.0	67.8	23.3	35.4	46.4	54.7	40.0	66.0	76.8	83.0	86.7	78.1	24.9	31.9	40.6	49.5	36.7
TOVA	61.8	73.6	78.6	84.0	74.5	16.4	32.2	42.1	48.6	34.8	71.6	76.8	79.1	81.9	77.3	21.6	29.5	35.1	42.2	32.1
+ OBCACHE-V	64.9	74.8	80.0	85.4	76.3	18.8	34.6	44.3	50.7	37.1	72.3	76.9	79.9	83.2	78.1	22.3	30.9	36.8	44.0	33.5
+ OBCACHE-K	65.3	75.0	80.2	85.7	76.5	20.3	36.9	46.1	52.5	39.0	73.0	78.0	81.1	83.1	78.8	23.7	30.5	37.5	44.9	34.1
+ OBCACHE-V&K	65.5	75.3	80.1	85.8	76.7	21.5	37.0	46.0	52.3	39.2	72.7	77.9	81.4	83.2	78.8	23.7	30.9	37.8	44.7	34.3
SnapKV	56.4	71.5	78.8	83.0	72.4	19.4	34.2	44.5	53.4	37.9	69.2	78.9	83.2	86.0	79.3	24.5	31.1	39.0	47.9	35.6
+ OBCACHE-V	55.8	72.2	77.6	82.0	71.9	18.6	32.1	41.7	50.7	35.8	68.5	78.0	83.7	86.6	79.2	24.3	30.3	37.7	46.3	34.7
+ OBCACHE-K	59.7	74.0	79.2	82.9	73.9	24.1	38.9	48.3	57.0	42.1	69.4	78.8	83.8	86.6	79.7	26.0	32.1	38.8	48.0	36.2
+ OBCACHE-V&K	59.0	73.6	78.9	82.9	73.6	22.4	38.5	48.8	58.0	41.9	69.7	79.3	83.6	86.5	79.8	26.2	31.9	39.3	48.2	36.4
AdaKV	58.9	75.4	82.4	86.0	75.7	21.5	39.1	50.5	60.7	43.0	72.2	81.9	86.1	87.5	81.9	26.9	38.2	52.6	64.3	45.5
+ OBCACHE-V	66.2	81.4	86.0	87.6	80.3	23.9	49.0	61.9	70.9	51.4	76.5	85.4	87.7	88.6	84.6	30.3	46.1	60.9	72.4	52.4
+ OBCACHE-K	69.2	82.6	86.3	88.3	81.6	29.2	52.0	65.5	73.2	55.0	78.1	86.1	88.0	88.5	85.2	31.7	47.5	63.6	74.8	54.4
+ OBCACHE-V&K	70.1	82.7	86.3	88.4	81.9	30.0	52.2	65.1	73.6	55.2	78.2	86.1	88.1	88.6	85.2	32.6	48.1	64.2	75.4	55.1

- **Integration into 4 Attn-only Baselines:** H₂O, TOVA, SnapKV, AdaKV.

- **Static Prefill Cache Eviction:** consistently improves RULER (4K, 32K) & LongBench performance under both query-aware & -agnostic settings, for both LLaMA-3.1-8B & Qwen-2.5-7B, at all compression rates.

- **Dynamic Decoding Eviction:** consistently reduces PG19 language modeling ppl.

- **Comparison of OBCache Scores:**

- value score: simplest, minimal overhead.
- key-aware scores: more output-sensitive, more accurate token saliency estimation.

Dynamic Decoding Eviction

