# REVIEW OF "CARTRIDGES: LIGHTWEIGHT AND GENERAL-PURPOSE LONG CONTEXT REPRESENTATIONS VIA SELF-STUDY" FOR ICLR 2025

Pavel Bushuyeu

October 2025

## 1 Summary

The paper introduces a method for representing long contexts more compactly by replacing the large runtime KV cache with a smaller, offline-trained CARTRIDGE. A CARTRIDGE is a learned key–value representation that acts as a persistent stand-in for the full KV cache normally constructed when the entire corpus is placed in context. Each CARTRIDGE is trained once for a specific document or corpus and can later be loaded at inference to reproduce the behavior of the model as if the full text were present. Training uses a SELF-STUDY process: the model generates synthetic question-answer dialogues about the corpus, and its next-token distributions are distilled into the smaller CARTRIDGE. The authors report that this approach matches ICL performance while using $38.6\times$ less memory and achieves $26.4\times$ higher throughput. They also demonstrate context extension beyond the model's native window (e.g., from 128K to 484K tokens on MTOB) and show that independently trained CARTRIDGES can be composed for multi-document queries. Experiments cover several long-document QA and comprehension benchmarks (Multi-Key Needle-in-a-Haystack, LONGHEALTH, MTOB, QASPER) with detailed ablations on factors such as prefix-tuning vs. LoRA, context-distillation vs. next-token training, and seed-prompt diversity.

## 2 Soundness

Score: 3/4

The technical framework of the paper, optimizing a small set of key-value vectors through prefix tuning and distillation, is conceptually sound and grounded in prior work (Li Liang 2021; Snell et al. 2022). The mathematical formulations of the distillation loss and the prefix parameterization are well-motivated. However, several aspects weaken the overall soundness:

1. Evaluation design - All CARTRIDGES are trained and tested in the same fixed corpus, which risks overfitting and inflating claims of 'general purpose' behavior. The absence of out-of-distribution or multi-domain evaluations makes it unclear whether the learned representations generalize.

2. Baseline scope - recent compression methods (e.g., Minicache 2024, PALU 2024) and architectural long-context alternatives (e.g., Mamba, Hyena, RWKV, xLSTM) are omitted. Without these, it is impossible to verify the claimed superiority.

3. Compute trade-offs - The paper acknowledges but underplays the cost of off-line training (2 to 4 orders of magnitude more FLOPs than ICL pre-fill). Claims of "cheaper serving" might therefore be misleading without further examples and analysis.

Overall, while the algorithmic idea is internally consistent and reproducible in principle, the empirical methodology leaves several gaps. The approach appears plausible but not fully validated as a viable alternative to long-context models.

# 3   Presentation

Score: 3/4

The paper is clearly written and well-structured. The motivation for CARTRIDGES is intuitive and the key ideas are presented in a logical manner. Figures effectively communicate the high-level intuition and empirical results.

However, the contextualization relative to prior work is weaker. The related work section lists many references but does not clearly articulate how CARTRIDGES differ in assumptions and limitations from KVPress or DuoAttention.

Overall, the presentation quality is solid, with strong visual design and coherent flow, but falls short of "excellent" due to limited contextual depth and a tendency to oversimplify some trade-offs in the main text.

# 4   Contribution

Score: 3/4

The paper addresses a highly relevant and timely problem and proposes a novel idea: replacing dynamic KV caches with smaller, offline-trained CARTRIDGES that can be reused across queries. This direction is elegant and draws a line between runtime attention mechanisms and offline-learned memory representations.

However, the overall contribution is somewhat limited by validation. Missing comparisons to leading KV-compression baselines (Lexico, Minicache, PALU) and long-context architectures (Mamba, Hyena, RWKV) weakens the impact. The offline costs might be substantial, which also diminishes the claim of practical efficiency.

Despite its limitations, the core idea is valuable. With stronger empirical grounding, this work could form the basis for a new class of learned memory systems and help clarify the boundary between general world knowledge encoded in the base model and corpus-specific knowledge.

# 5    Strengths

- **Sound idea:** The central idea of replacing dynamically built KV caches with offline-trained, reusable CARTRIDGES is technically elegant and well-grounded in prior work on prefix-tuning and distillation. It introduces a clear separation between the runtime inference mechanism and an offline learned memory representation, which is both original and practically motivated.

- **Strong empirical gains:** The paper reports substantial reductions in memory usage (up to 38.6$\times$) and large improvements in throughput (up to 26.4$\times$) while maintaining comparable quality to full in-context learning across multiple benchmarks. These results convincingly demonstrate that the method achieves a meaningful trade-off between efficiency and performance.

- **Comprehensive ablation analysis:** The experimental section is thorough. The authors systematically evaluate several dimensions of their approach, including prefix-tuning versus LoRA parameterization, next-token prediction versus distillation objectives, and the impact of synthetic prompt diversity. This level of rigor demonstrates a strong understanding of how each design choice affects performance, lending credibility to the empirical findings.

- **Practical relevance:** The work directly targets a pressing bottleneck in large language model deployment. By reframing the challenge as one of amortized learning rather than purely runtime optimization, the paper proposes a new direction of thought about memory management in LLMs.

- **Clarity and readability:** The paper is well-written, logically organized, and visually clear. Figures effectively communicate both the motivation and the results. The introduction and method sections have a good balance between intuition and formalism, making the paper accessible to a broad audience.

- **Forward-looking potential:** Beyond immediate results, the proposed paradigm opens new research territory at the intersection of learned retrieval, knowledge distillation, and context compression. The notion of maintaining persistent, corpus-specific CARTRIDGES invites exploration into questions of how models separate general world knowledge from corpus-specific adaptation—a promising and underexplored direction for the larger community.

# 6 Weaknesses

- **Missing baselines:** The paper omits several major methods in both KV-cache compression and long-context modeling. Notably absent are Lexico (2024), Minicache (2024), PALU (2024), and KVPress, as well as architectural alternatives such as Mamba, Hyena, RWKV, and xLSTM. Without these, the empirical claim of state-of-the-art efficiency remains open.

- **Narrow evaluation scope:** Although the Appendix lists several potential application types for CARTRIDGES (e.g., summarization, retrieval, and multi-document reasoning), the main experiments only evaluate single-document QA and translation tasks. The additional use cases are described but not empirically demonstrated or analyzed. Consequently, the claim that CARTRIDGES provide "general-purpose" long-context representations remains speculative.

- **Compute cost distortion:** Although the paper reports large gains in inference efficiency, it underplays the offline training costs, which the authors acknowledge can be two to four orders of magnitude higher in FLOPs than a standard ICL prefill. This imbalance undermines the claim of overall cost efficiency.

- **Synthetic-data dependency:** The SELF-STUDY approach relies on automatically generated QA pairs whose diversity and realism are unclear. The paper does not specify the prompting strategy, model choice, or quantity of synthetic data used. This omission makes it difficult to evaluate how robust or representative the learned CARTRIDGES truly are.

- **Overfitting risk:** Each CARTRIDGE is trained and tested on the same document or corpus, raising concerns that the method captures document-specific memorization rather than generalizable reasoning behavior. Evaluations on out-of-domain corpora or held-out queries would help clarify this limitation.

- **Composition and reproducibility issues:** The procedure for combining multiple CARTRIDGES at inference time is only briefly described and lacks technical detail. Without clearer documentation or released scripts, the reproducibility of these compositional results is questionable.

# 7 Questions

- **Synthetic data generation:** You mention that the SELF-STUDY dialogues were generated using Claude. Could you provide more detail on this process—specifically the prompt templates, the number of examples per document, and whether diversity or filtering mechanisms were used?

It would also be helpful to know how performance varies with the amount or source of synthetic data.

- **Compute and efficiency trade-offs:** The paper acknowledges that training CARTRIDGES can require two to four orders of magnitude more FLOPs than a single ICL prefill. Could you clarify how this compares to standard fine-tuning in terms of total GPU hours? At what scale of query reuse does the method become cost-effective?

- **Evaluation breadth:** The appendix lists several potential applications (e.g., summarization, retrieval, reasoning), yet the main results only cover QA and translation tasks. Were there any qualitative checks on these other use cases, even if not reported? If not, could you comment on any limitations that prevented such evaluations?

- **Baseline selection:** Could you elaborate on the rationale for not comparing against compression baselines such as Lexico, Minicache or PALU?

- **Composition mechanism:** The paper briefly notes that multiple CARTRIDGES can be composed at inference time, but the mechanism is unclear. Are there limitations in compositional depth or order sensitivity?

- **Composition within the same corpus:** Have you explored using multiple CARTRIDGES trained on different partitions of the same corpus, either to improve coverage or enable parallelization? If so, how does this affect performance and consistency across queries?

# 8 Flag for Ethics Review

No ethics review needed

# 9 Rating

Score: 6

# 10 Confidence

Score: 2-3