Context Adaptive Memory-Efficient LLM Inference for Edge Multi-Agent Systems

Extended Abstract

Hamza Mohammed Samsung Research America Mountain View, CA, USA h.mohammed@samsung.com Hang Yin Samsung Research America Mountain View, CA, USA hang.y@samsung.com Sai Chand Boyapati Samsung Research America Mountain View, CA, USA s.boyapati@samsung.com

ABSTRACT

Large Language Models (LLMs) excel at multi-document QA, summarization, code generation, and other language-intensive tasks, yet they demand substantial memory resources for storing keyvalue (KV) caches and processing attention in long-context scenarios. These requirements often prohibit on-device or edge deployments in *multi-agent systems (MAS)*, where multiple agents share or update contextual information and need efficient inference pipelines. We present CASK (Context-Adaptive Sparse Key-value), an *inference-time* strategy that reduces memory usage while preserving strong performance on extended contexts. CASK addresses this challenge with two complementary mechanisms: a dynamic sparse attention module-a lightweight, meta-learned component-that identifies the most relevant context tokens, and an adaptive KV-cache compression technique that dynamically quantizes and prunes less critical key-value pairs based on usage frequency and recency. These innovations enable near-lossless performance on long-context tasks while cutting memory usage by up to 40% and boosting inference speed by as much as 20%. Evaluations on LongBench [2] and multi-agent benchmarks show that CASK maintains over 95% of baseline accuracy while allowing more agents or extended histories under tight GPU budgets. Integration into a vision-language agent for collaborative, multimodal contexts underscores its practicality for resource-constrained LLM deployments in MAS.

KEYWORDS

Multi-Agent Systems; Natural Language Generation; Dimensionality Reduction And Manifold Learning; Heuristic Function Construction; Caching And Paging Algorithms; Multi-Agent Learning; Computer Vision; Reinforcement Learning; Multi-Agent Reinforcement Learning; Active Learning; Neural Networks; Supervised Learning; Online Learning Algorithms

ACM Reference Format:

Hamza Mohammed, Hang Yin, and Sai Chand Boyapati. 2025. Context Adaptive Memory-Efficient LLM Inference for Edge Multi-Agent Systems: Extended Abstract. In Proc. of the 24th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2025), Detroit, Michigan, USA, May 19 – 23, 2025, IFAAMAS, 3 pages.

This work is licensed under a Creative Commons Attribution International 4.0 License.

1 INTRODUCTION

Large Language Models (LLMs) have dramatically advanced natural language processing, enabling high performance in few-shot learning, code synthesis, and long-form text generation [3, 4, 8, 18]. However, these models are challenging to deploy in *multi-agent systems (MAS)* with limited computational resources [14, 19, 29]. Each agent may need to maintain an extensive context history or share contextual knowledge with others, magnifying the problem of storing large key-value (KV) caches and executing the quadratic attention mechanism in standard Transformers [3, 26]. This restricts real-time or on-device inference in settings such as IoT networks, robotic swarms, or other privacy and latency-critical domains.



Figure 1: The Vision Language Agent example for MAS, a vision-language model integrated with CASK. The Mask Generation Module (MGM) periodically outputs sparse attention masks and triggers KV-cache compression.

We introduce **CASK (Context-Adaptive Sparse Key-value)** an *inference-time* approach that substantially cuts memory usage while retaining the core strengths of LLM-based systems in extended-context tasks. Unlike many existing techniques that require re-training or architecture modifications [5, 6, 11, 21], CASK operates on top of *pre-trained* LLMs. It dynamically generates sparse attention patterns and aggressively compresses the KV-cache to accommodate long contexts under strict memory constraints. Critically, it preserves coherent multi-agent interactions by focusing on crucial cross-references among agents' dialogue or sensor data.

Proc. of the 24th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2025), Y. Vorobeychik, S. Das, A. Nowé (eds.), May 19 – 23, 2025, Detroit, Michigan, USA. © 2025 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org).

2 APPROACH

CASK introduces two main components (Figure 1). The first is Dynamic Sparse Attention, which employs a Mask Generation Module (MGM), a small vision transformer [10] trained to derive sparse binary masks that eliminate lower-importance interactions via meta reinforcement learning from reconstruction [9, 27]. We apply these sparse masks to the attention logits, allowing the model's existing attention pipeline to leverage sparse operations without altering its underlying architecture. This effectively reduces the attention window, lowering memory and compute overhead. To handle contexts longer than the model's original training window, we include a positional interpolation scheme (e.g., [16, 17]) that smoothly extends the receptive field without requiring major retraining.

The second component is Adaptive KV-Cache Compression. CASK extends post-training weight-quantization methods [13, 22, 25] to dynamic KV-Cache quantization and pruning. By tracking recency, access frequency, attention allocation, and gradient attribution, it computes saliency metrics for each key-value pair. Guided by two thresholds—validated via runtime reconstruction loss—pairs with moderately low salience are dynamically quantized to a lower bit width, while those below a stricter threshold are pruned. This reduces cache size as contexts accumulate and is especially beneficial in MAS settings, where large histories often contain irrelevant or outdated content.

3 EXPERIMENTAL RESULTS

We implement a MAS with Vision Language Agents by integrating CASK onto a LLaMA 3.2 90B instruct model [8] as shown in Figure 1, wherein multiple agents communicate both text and visual embeddings, sharing a global context. To demonstrating CASK's stable long-context processing, critical for MAS-focused tasks, we evaluate on *LongBench* [2].

Table 1: Average Performance on LongBench [2].

Avg. Score (LongBench)	Memory Usage (GB)	Relative Inference Time
63.2	72	1.00x
44.0	350*	0.85x
56.4	63	0.92x
48.9	88	1.15x
52.6	43	0.78x
57.2	40	0.88x
62.8	68.5	0.82x
54.6	68.1	0.82x
37.1	40.8	0.88x
62.7	67.4	0.81x
61.5	44.5	0.77x
	Avg. Score (LongBench) 63.2 44.0 56.4 48.9 52.6 57.2 62.8 54.6 37.1 62.7 61.5	Avg. Memory Score Usage (LongBench) (GB) 63.2 72 44.0 350* 56.4 63 48.9 88 52.6 43 57.2 40 62.8 68.5 54.6 68.1 37.1 40.8 62.7 67.4 61.5 44.5

*Estimated usage based on model size.

4 **DISCUSSION**

Tables 1 and 2 compare our proposed CASK method against state-ofthe-art LLM test-time optimizations and the full-attention baseline. Overall, CASK achieves a strong balance of reduced memory usage, improved inference speed, and minimal drops in accuracy relative to the baseline. Specifically, in Table 1, CASK retains 61.5 average score

Table 2: Task-specific performance (LLaMA-3.2-90B-128k [8] on LongBench [2]).

Technique	S-Doc QA	M-Doc QA	Summ.	Few-shot	Code	Synth.
Full Attention	51.2	60.5	49.6	78.2	76.4	63.0
SEA	47.1	56.7	45.6	74.5	72.6	59.3
BigBird	50.6	56.1	48.6	66.7	68.3	37.1
StreamingLLMs	26.6	40.5	45.6	50.7	55.0	4.4
SampleAttention	49.5	60.4	49.4	78.1	76.3	62.9
CASK (Ours)	49.4	58.9	45.6	76.3	76.1	62.8

while cutting memory usage by 38% relative to the base LLaMA model [8], accelerating inference by around 23%. This improvement is consistent across input lengths up to 128k tokens, demonstrating stable long-context processing.

Looking at the breakdown in Table 2, CASK remains close to full attention on tasks requiring more intricate or broad context dependencies (e.g., code generation, few-shot, or synthetic tasks). Its chunked/compressed attention mechanism captures long-range interactions without incurring the quadratic memory cost of standard attention. By contrast, for narrower tasks (e.g., single-document QA, summarization), the localized context dependencies are often well-served by more specialized mechanisms, which can explain why certain methods such as BigBird (block-sparse) [26] or StreamingLLMs (step-by-step) [23] may outperform CASK in those cases. These approaches aggressively reduce attention overhead for short or more localized sequences, often matching or beating CASK's memory usage.

However, on tasks demanding broader or more complex attention patterns (e.g., multi-document QA, code generation, and synthetic benchmarks), simpler or more aggressively sparse methods can lose critical long-range interactions. Here, CASK's chunk-based strategy offers more flexibility in capturing essential dependencies, allowing it to surpass the more specialized methods. Other techniques like SEA [12] or SampleAttention [30] may track closely to CASK in some benchmarks, yet they typically incur approximately 30% higher memory usage. Hence, CASK emerges as an attractive trade-off for real-world deployments where memory constraints and inference speed are key concerns, while still delivering strong overall performance across a variety of tasks.

5 CONCLUSION

CASK addresses a central challenge in deploying LLMs for *multiagent systems* with long or evolving contexts under strict resource budgets. By leveraging dynamic sparse attention and importancebased KV-cache compression, it (i) maintains robust performance (over 95% of baseline accuracy) on diverse tasks requiring significant context windows, (ii) saves ~40% of GPU memory and reduces latency (allowing more agents or deeper contexts without out-ofmemory errors), and (iii) integrates into existing LLM inference stacks with minimal changes—requiring no specialized retraining or hardware. Future work includes exploring automated threshold tuning for compression, investigating cross-agent attention sharing, and tailoring the method to specialized edge hardware. Overall, CASK provides a practical solution for memory-efficient LLM inference in resource-constrained, collaboration-focused environments.

REFERENCES

- [1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond. arXiv preprint arXiv:2308.12966 (2023).
- [2] Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. LongBench: A Bilingual, Multitask Benchmark for Long Context Understanding. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, Bangkok, Thailand, 3119–3137. https://doi.org/10.18653/v1/2024.acllong.172
- [3] Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The Long-Document Transformer. (2020). arXiv:2004.05150 [cs.CL] https://arxiv.org/abs/ 2004.05150
- [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Proceedings of the 34th International Conference on Neural Information Processing Systems (Vancouver, BC, Canada) (NIPS '20). Curran Associates Inc., Red Hook, NY, USA, Article 159, 25 pages.
- [5] Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. 2023. Extending Context Window of Large Language Models via Positional Interpolation. arXiv:2306.15595 [cs.CL] https://arxiv.org/abs/2306.15595
- [6] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. 2019. Generating Long Sequences with Sparse Transformers. arXiv:1904.10509 [cs.LG] https: //arxiv.org/abs/1904.10509
- [7] Krzysztof Marcin Choromanski, Valerii Likhosherstov, David Martin Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. 2021. Rethinking Attention with Performers. In accepted to ICLR 2021 (oral presentation). https://arxiv.org/abs/2009.14794
- [8] Abhimanyu Dubey et al. 2024. The Llama 3 Herd of Models. arXiv preprint arXiv:2407.21783 (2024). https://arxiv.org/abs/2407.21783
- [9] Haotian Fu, Hongyao Tang, Jianye Hao, Chen Chen, Xidong Feng, Dong Li, and Wulong Liu. 2021. Towards Effective Context for Meta-Reinforcement Learning: an Approach based on Contrastive Learning. In *Thirty-Fifth AAAI Conference* on Artificial Intelligence. AAAI Press, 7457–7465. https://ojs.aaai.org/index.php/ AAAI/article/view/16914
- [10] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. 2023. Segment Anything. In 2023 IEEE/CVF International Conference on Computer Vision (ICCV). 3992–4003. https://doi.org/10.1109/ ICCV51070.2023.00371
- [11] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. 2020. Reformer: The Efficient Transformer. In International Conference on Learning Representations. https: //openreview.net/forum?id=rkgNKkHtvB
- [12] Heejun Lee, Jina Kim, Jeffrey Willette, and Sung Ju Hwang. 2024. SEA: Sparse Linear Attention with Estimated Attention Mask. In *The Twelfth International Conference on Learning Representations*. https://openreview.net/forum?id=JbcwfmYrob
- [13] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. 2024. AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration. In *MLSys*.
- [14] Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative Agents: Interactive Simulacra of Human Behavior. In Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (San Francisco, CA, USA) (UIST '23). Association for Computing Machinery, New York, NY, USA, Article 2, 22 pages. https: //doi.org/10.1145/3586183.3606763
- [15] Andrew Peng, Michael Wu, John Allard, Logan Kilpatrick, and Steven Heidel. 2023. GPT-3.5 Turbo Fine-Tuning and API Updates. https://openai.com/index/gpt-3-5-turbo-fine-tuning-and-api-updates.
- [16] Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2024. YaRN: Efficient Context Window Extension of Large Language Models. In The Twelfth

 $\label{eq:linear} International \ Conference \ on \ Learning \ Representations. \ https://openreview.net/forum?id=wHBfxhZu1u$

- [17] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. RoFormer: Enhanced transformer with Rotary Position Embedding. *Neurocomput.* 568, C (Feb. 2024), 12. https://doi.org/10.1016/j.neucom.2023.127063
- [18] H. Touvron, T. Lavril, G. Izacard, and et al. 2023. LLaMA: Open and efficient
- foundation language models. arXiv preprint arXiv:2302.13971 (2023).
 [19] Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Dongdong Zhang, and Furu Wei. 2024. DeepNet: Scaling Transformers to 1,000 Layers. *IEEE Trans. Pattern Anal. Mach. Intell.* 46, 10 (Oct. 2024), 6761–6774. https://doi.org/10.1109/ TPAMI.2024.3386927
- [20] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution. arXiv preprint arXiv:2409.12191 (2024).
- [21] Yingsheng Wu, Yuxuan Gu, Xiaocheng Feng, Weihong Zhong, Dongliang Xu, Qing Yang, Hongtao Liu, and Bing Qin. 2024. Extending Context Window of Large Language Models from a Distributional Perspective. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 7288–7301. https://doi.org/10.18653/v1/2024.emnlpmain.414
- [22] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. SmoothQuant: Accurate and Efficient Post-Training Quantization for Large Language Models. In Proceedings of the 40th International Conference on Machine Learning.
- [23] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient Streaming Language Models with Attention Sinks. In *The Twelfth International Conference on Learning Representations*. https://openreview.net/forum? id=NG7sS51zVF
- [24] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. Qwen2 Technical Report. arXiv:2407.10671 [cs.CL] https://arxiv.org/abs/2407.10671
- [25] Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong He. 2022. ZeroQuant: efficient and affordable post-training quantization for large-scale transformers. In Proceedings of the 36th International Conference on Neural Information Processing Systems (New Orleans, LA, USA) (NIPS '22). Curran Associates Inc., Red Hook, NY, USA, Article 1970, 16 pages.
- [26] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. Big bird: Transformers for longer sequences. Advances in Neural Information Processing Systems 33 (2020).
- [27] Zhejun Zhang, Alexander Liniger, Dengxin Dai, Fisher Yu, and Luc Van Gool. 2021. End-to-End Urban Driving by Imitating a Reinforcement Learning Coach. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV).
- [28] Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark Barrett, Zhangyang Wang, and Beidi Chen. 2023. H2O: heavy-hitter oracle for efficient generative inference of large language models. In Proceedings of the 37th International Conference on Neural Information Processing Systems (New Orleans, LA, USA) (NIPS '23). Curran Associates Inc., Red Hook, NY, USA, Article 1506, 50 pages.
- [29] Zhenyu Zhou, Yi Tay, Ramesh Nallapati, Bhaskar Mitra, Zhicheng Xiao, Hao Cheng, Xiangru Xiang, Joseph P Sim, Harish Swaminathan, Nam D Tran, et al. 2022. Efficient language modeling with sparse all-MLP. arXiv preprint arXiv:2203.06850 (2022).
- [30] Qianchao Zhu, Jiangfei Duan, Chang Chen, Siran Liu, Xiuhong Li, Guanyu Feng, Xin Lv, Huanqi Cao, Xiao Chuanfu, Xingcheng Zhang, Dahua Lin, and Chao Yang. 2024. SampleAttention: Near-Lossless Acceleration of Long Context LLM Inference with Adaptive Structured Sparse Attention. *CoRR* abs/2406.15486 (2024). https://doi.org/10.48550/ARXIV.2406.15486