Towards Efficient Large Language Model Serving: A Survey on System-Aware KV Cache Optimization

Jiantong Jiang¹, Peiyu Yang¹, Rui Zhang², Feng Liu¹

¹The University of Melbourne, ²Huazhong University of Science and Technology {jiantong.jiang, peiyu.yang}@unimelb.edu.au rayteam@yeah.net, fengliu.ml@gmail.com

Abstract

Despite the rapid advancements of large language models (LLMs), LLM serving systems remain memory-intensive and costly. The keyvalue (KV) cache, which stores KV tensors during autoregressive decoding, is crucial for enabling low-latency, high-throughput LLM inference serving. In this survey, we focus on system-aware KV infrastructure for serving LLMs (abbreviated as sKis). We revisit recent work from a system behavior perspective, organizing existing efforts into three dimensions: execution and scheduling (temporal), placement and migration (spatial), and representation and retention (structural). Furthermore, we analyze synergies across behaviors and their links to optimization objectives, highlighting future opportunities. Our work systematizes a rapidly evolving area, providing a foundation for understanding and innovating KV cache designs in modern LLM serving infrastructure.

1 Introduction

Large language models (LLMs) have showcased exceptional abilities across diverse applications (Zhao et al., 2023), with notable examples like GPT (Radford et al., 2018, 2019; Brown et al., 2020; OpenAI, 2023), LLaMA (Touvron et al., 2023a,b), and OPT (Zhang et al., 2022). These models excel at large-scale high-quality language understanding and generation, powered by the Transformer architecture (Vaswani et al., 2017), which efficiently captures long-range dependencies via self-attention.

Despite their success, serving LLMs efficiently remains non-trivial (Li et al., 2024a). Transformer-based LLMs generate tokens autoregressively, with each token conditioned on all previous ones. To avoid redundant compute, serving systems adopt a *key-value (KV) cache* (Pope et al., 2023) to store intermediate KV tensors of the generated tokens. Yet, as prompt and output length grow, the KV cache can reach 1–2 million tokens (Pichai and

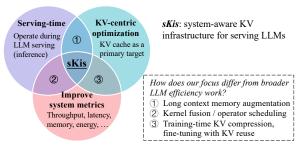


Figure 1: Positioning of the survey scope ("sKis").

Hassabis, 2024; Ding et al., 2024), creating memory bottlenecks and highlighting the critical role of KV cache optimization. Consequently, a growing body of KV-centric techniques has emerged, yielding memory savings and efficiency gains in throughput and latency (Li et al., 2024b).

To this end, we argue that it deserves a deep investigation of system-aware, serving-time, KVcentric optimization methods, as shown in Fig. 1, which we call this scope sKis. We adopt a systemoriented taxonomy to offer a comprehensive understanding of sKis, categorizing methods along three fundamental axes of system behaviors, as shown in Fig. 2: (i) execution and scheduling focuses on the temporal control of when KV data is accessed, computed, or scheduled (cf. § 3); (ii) placement and migration captures the spatial decisions of where KV data is placed in the memory hierarchy or migrated across devices (cf. § 4); and (iii) representation and retention concerns the structural treatment of how KV data is compressed or managed (cf. § 5). We further analyze cross-behavior synergies and behavior-objective effects to reveal overlooked regions and open challenges (cf. § 6).

Prior surveys span efficient LLM inference and serving (Zhou et al., 2024b; Yuan et al., 2024; Miao et al., 2023; Li et al., 2024a; Wang et al., 2024a), with KV cache discussed only as a minor component. Another line specifically reviews KV cache methods, organizing them by lifecycle stages (Shi et al., 2024) or by abstraction level (Li et al., 2024b).

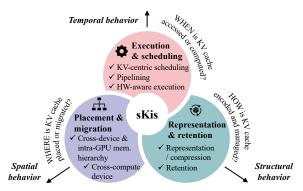


Figure 2: Taxonomy of the survey that covers temporal, spatial, and structural dimensions.

Instead, this survey focuses exclusively on sKis and distinguishes itself by offering a system-oriented perspective and a deeper understanding. Comparisons with related surveys are in App. B. To the best of our knowledge, we are the first to frame KV cache optimization as a temporal-spatial-structural behavior space, which enables comparable synthesis and actionable guidance for this area. We hope our work informs future research in LLM acceleration, especially that centered on KV cache design.

2 Foundations and Taxonomy

Transformer-based LLMs. LLMs are built from stacked Transformer blocks, each with multi-head self-attention (MHSA) and feed-forward network (FFN). These blocks are sequential, where the output of one block serves as the input to the next. For the i-th attention head, MHSA applies learned projections W^{Q_i} , W^{K_i} , and W^{V_i} to the input features X to get queries $Q_i = XW^{Q_i}$, keys $K_i = XW^{K_i}$, and values $V_i = XW^{V_i}$. The self-attention operation then generates the head output Z_i , and the outputs from all heads are concatenated and projected by W^O to form the final attention output:

$$Z_i = \operatorname{softmax}(\frac{Q_i K_i^{\top}}{\sqrt{d_k}}) V_i, Z = \operatorname{concat}(Z_1, ..., Z_h) W^O.$$

Following this, the output of MHSA is fed into the FFN module, which applies two linear transformations with a nonlinear activation (e.g., ReLU):

$$FFN(x) = ReLU(xW_1 + b_1)W_2 + b_2,$$

where W_1 , W_2 , b_1 , and b_2 are learnable parameters. These modules together enable contextualized autoregressive modeling in LLMs.

LLM Inference and KV Cache. LLMs generate tokens in an autoregressive manner, as shown in Fig. 3. In each generation step, the model takes as input the input tokens and previously generated

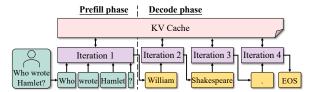


Figure 3: Prefill and decode phases of LLM inference.

tokens, and generates the next token. This process is divided into two phases: the *prefill* phase starts with the initial input tokens and generates the first output token, while the *decode* phase generates new tokens autoregressively. Due to the quadratic time complexity of self-attention, repeatedly computing attention across tokens is computationally expensive. To address this, *key-value* (*KV*) cache is used to store the previously computed intermediate KV tensors, as shown in Fig. 3. At each new step, the model efficiently reuses these cached tensors without recomputing attention over the entire sequence. **Scope and Taxonomy.** This survey investigates recent advances in the sKis scope shown in Fig. 1.

sKis denotes system-aware KV infrastructure for serving LLMs. A method belongs to sKis if it: (i) operates during serving (inference), (ii) centers on KV caches as the primary optimization target, and (iii) aims to improve system metrics without retraining the base LLM's weights or modifying its Transformer architecture.

This survey organizes literature on sKis via a system view, as shown in Fig. 2. Unlike task- or stage-based classifications, our taxonomy is grounded in low-level system behaviors, and we offer further details in App. A. However, similar to how a modern OS includes components for scheduling, memory, and I/O, LLM serving systems often involve techniques spanning various aspects. Thus, a single paper may naturally touch on several categories. For clarity and focus, we highlight one or two primary categories per work based on its main contributions. Minor associations are not elaborated, and we refer to App. C for details. We summarize the methods in Fig. 4 and the key findings in App. D.

3 KV Execution and Scheduling

This section captures the temporal behaviors of KV cache usage, including how cache entries are scheduled and executed efficiently at runtime.

3.1 KV-centric Scheduling

While scheduling is a long-studied system problem, KV-centric scheduling (KVS) methods explic-

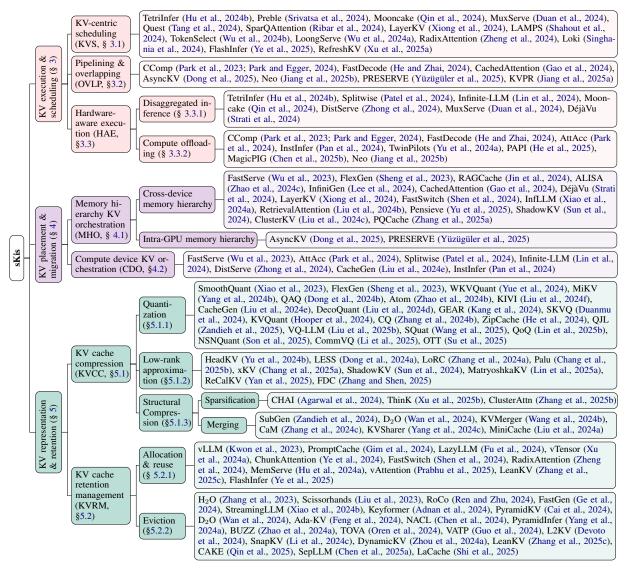


Figure 4: Taxonomy of sKis and associated methods. Each method is annotated with its primary contributions for conciseness. Minor category associations are omitted here and listed in Appx. C, Tab. 5.

itly integrate KV characteristics into runtime decisions. At the request level, some methods adopt KV usage-aware scheduling to balance resource load and reduce contention (Hu et al., 2024b; Duan et al., 2024; Xiong et al., 2024; Shahout et al., 2024; Wu et al., 2024a). For example, TetriInfer (Hu et al., 2024b) prioritizes requests using predicted KV usage to mitigate prefill-decode interference, while LoongServe (Wu et al., 2024a) integrates request dispatching, KV placement, batching, and parallelism scaling into a unified scheduling system guided by runtime profiling and cost modeling. Some methods prioritize requests with high KV reuse potential (Srivatsa et al., 2024; Zheng et al., 2024). Mooncake (Qin et al., 2024) considers prefix reuse and hotspot balance to select prefill-decode instance pairs for requests.

At finer granularity, methods schedule attention

at the **token level**, selecting KV based on attention contributions that are either query-aware (Tang et al., 2024; Ribar et al., 2024; Singhania et al., 2024) or query-agnostic (Wu et al., 2024b), or via periodic refresh that alternates full-context and subset attention (Xu et al., 2025a). At the **kernel level**, FlashInfer (Ye et al., 2025) schedules attention workloads across CUDA thread blocks based on query and KV lengths to reduce GPU idleness.

3.2 Pipelining and Overlapping

Pipelining and overlapping (*OVLP*) methods hide latency by concurrently executing KV-related compute, communication, and I/O, which are key to reducing idle time and improving serving efficiency. Though often embedded in broader systems, Tab. 1 highlights methods where OVLP forms the core technical contribution. Some distribute compute

Table 1: Summary of representative pipelining and overlapping methods.

Method	Overlapped operations (with device)	Туре
FastDecode (He and Zhai, 2024) Neo (Jiang et al., 2025b) CComp (Park and Egger, 2024) CachedAttention (Gao et al., 2024) KVPR (Jiang et al., 2025a) AsyncKV (Dong et al., 2025) PRESERVE (Yüzügüler et al., 2025)	CPU R-part compute \leftrightarrow GPU S-part compute CPU attention compute \leftrightarrow GPU linear ops CPU MHSA compute \leftrightarrow FFN data transfer (CPU \rightarrow GPU) KV load/store (CPU \leftrightarrow GPU) \leftrightarrow GPU compute GPU KV recompute \leftrightarrow KV transfer (CPU \leftrightarrow GPU) GPU KV prefetch (HBM \rightarrow L2) \leftrightarrow GPU attention compute GPU KV prefetch (HBM \rightarrow L2) \leftrightarrow GPU collective communication	Comp-Comp Comp-I/O Comp-I/O Comp-I/O Comp-I/O I/O-Comm

across devices, enabling pipelining between CPU and GPU (He and Zhai, 2024; Jiang et al., 2025b). Others overlap compute with host-device KV transfer or GPU memory I/O (Park and Egger, 2024; Gao et al., 2024; Dong et al., 2025; Jiang et al., 2025a). Recently, Yüzügüler et al. (2025) explored I/O-communication overlap by overlapping KV prefetch with inter-GPU collective communication.

3.3 Hardware-aware Execution

This section focuses on hardware-aware execution (*HAE*) methods that adapt KV-related operations to the underlying heterogeneous hardware.

3.3.1 Disaggregated Inference

Disaggregated inference separates heterogeneous computation in LLM inference and maps them to distinct hardware resources to reduce interference and improve utilization. Infinite-LLM (Lin et al., 2024) adopts this principle at the operator level, disaggregating attention across distributed instances.

Four concurrent papers propose to separate the prefill and decode phases due to their inherent asymmetry: prefill is compute-intensive, whereas decode is memory-bound (Hu et al., 2024b; Patel et al., 2024; Zhong et al., 2024; Strati et al., 2024). They assign each phase to a separate set of GPU instances, with KV caches reused across phases via transfer or memory sharing. Mooncake (Qin et al., 2024) takes a softer form, scheduling requests with different KV access patterns and workloads to different GPUs. By comparison, MuxServe (Duan et al., 2024) also decouples the phases but colocates them within a single GPU unit via SM partitioning.

3.3.2 Compute Offloading

Compute offloading relocates partial compute to auxiliary devices to reduce GPU bottlenecks, utilizing hardware heterogeneity and workload features. They often follow a compute-near-cache principle, placing compute near KV caches for better locality. Early work on **CPU offloading** assigns partial decoding to the CPU (Park et al., 2023), later ex-

tended to dynamic GPU-CPU division (Park and Egger, 2024). TwinPilots (Yu et al., 2024a) divides transformer layers into operations and solves a linear program to assign them to GPU or CPU. Fast-Decode (He and Zhai, 2024) and Neo (Jiang et al., 2025b) offload attention and KV caches, using costaware hardware selection and a load-aware scheduler, respectively. MagicPIG (Chen et al., 2025b) also offloads attention and its retrieval tasks.

Beyond CPUs, several methods offload compute to **alternative devices**. For instance, InstInfer (Pan et al., 2024) offloads memory-intensive attention to the computational storage drive (CSD), exploiting the high bandwidth of flash chips. Some works also explore processing-in-memory (PIM)-based offloading (He et al., 2025; Park et al., 2024). These approaches expand the compute offloading space towards broader device heterogeneity.

4 KV Placement and Migration

This section focuses on the spatial behaviors of how KV caches are placed and migrated across memory hierarchies and between compute devices. Figure 5 visualizes the architecture and transfer paths.

4.1 Memory Hierarchy KV Orchestration

To scale under memory limits, we survey memory hierarchy KV orchestration (*MHO*) methods that distribute KV caches across memory hierarchies.

One common direction migrates KV entries between **cross-device memory hierarchies** to reduce GPU memory pressure. These methods manage KV caches across fast but limited GPU HBM, and larger but slower alternatives like CPU DRAM or SSD. Some works focus on **importance-aware** migration, where only critical KV entries are kept on GPU. ALISA (Zhao et al., 2024c) and Pensieve (Yu et al., 2025) estimate token importance online to offload less useful ones and recover them on demand. InfiniGen (Lee et al., 2024) offloads all KV entries to the CPU and prefetches critical ones via speculated attention scores. ShadowKV (Sun et al., 2024) likewise offloads all value caches and

fetches selected chunks guided by key-side landmarks. InfLLM (Xiao et al., 2024a) retains only local and frequently used KV blocks on the GPU. ClusterKV (Liu et al., 2024c) and RetrievalAttention (Liu et al., 2024b) also offload the majority of KV caches. They then retrieve semantically relevant KV entries or clusters to reload. Similarly, PQCache (Zhang et al., 2025a) compresses all KV entries on the CPU and retrieves relevant ones using approximate attention scores. Other methods adopt system cost-driven strategies, guiding KV placement and migration via cost models or runtime trade-offs. FlexGen (Sheng et al., 2023) places KV caches across GPU, CPU, and disk by formulating a cost model that optimizes throughput under bandwidth and latency constraints. In contrast, many systems make runtime decisions of KV offloading, prefetching, or recovery based on activity status (Wu et al., 2023; Strati et al., 2024), access patterns (Jin et al., 2024; Xiong et al., 2024), look-ahead hints (Gao et al., 2024), or scheduling events (Shen et al., 2024).

Another line of MHO methods targets the **intra-GPU memory hierarchy**. These methods migrates KV entries between on-chip L1/L2 caches and off-chip HBM to hide latency. Dong et al. (2025) asynchronously prefetched the upcoming KV blocks from HBM to the L2 cache during the current compute cycles. Similarly, PRESERVE (Yüzügüler et al., 2025) prefetches KV caches during collective communication, and inserts such operations selectively through graph-level optimization to avoid cache pollution.

4.2 Compute Device KV Orchestration

Unlike hierarchical memory, compute device KV orchestration (CDO) places and moves KV across compute-capable devices like GPUs, CPUs, and storage-attached processors, to enable distributed or heterogeneous serving. FastServe (Wu et al., 2023) divides KV caches across GPUs based on parallelism, and transfers intermediate results via peer-to-peer direct memory access (DMA). Other methods split the computations and distribute both computations and KV caches across GPUs (Patel et al., 2024; Lee et al., 2024; Zhong et al., 2024). For instance, DistServe (Zhong et al., 2024) proposes placement schemes for prefill and decode phases across both high and low node-affinity clusters, and uses a pull-based scheme where decode GPUs fetch KV entries as needed from prefill GPUs acting as buffers. Beyond tightly coupled

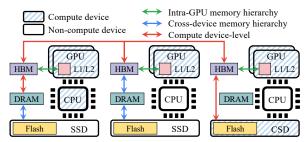


Figure 5: Illustration of KV cache placement and migration across memory hierarchies and compute devices.

GPU clusters, CacheGen (Liu et al., 2024e) targets remote KV cache transmission in distributed networked setups. It reduces network delay by compressing KV tensors into bitstreams and adaptively streaming them based on runtime bandwidth.

Extending this idea beyond GPU clusters, AttAcc (Park et al., 2024) and InstInfer (Pan et al., 2024) offload attention to PIMs or CSDs and support KV cache migration through optimized DMA.

5 KV Representation and Retention

This section focuses on structural system behaviors of KV cache representation and retention.

5.1 KV Cache Compression

With KV caches as the dominant memory bottleneck, KV cache compression (KVCC) is now a central research thrust.

5.1.1 KV Cache Quantization

Quantization compresses floating-point tensors into lower-precision formats. Early works enable 8-bit and 4-bit KV quantization with minimal quality loss (Xiao et al., 2023; Sheng et al., 2023). Some methods further adopt mixed-precision strategies, assigning higher precision to more critical KV entries based on certain criteria. We compare existing methods in Tab. 2 and present key insights here.

One insight is the **asymmetric KV quantization** due to their distinct outlier patterns and quantization sensitivities (Dong et al., 2024b; Liu et al., 2024f; Wang et al., 2025; Lin et al., 2025b), commonly quantizing keys per-channel and values pertoken (Liu et al., 2024f; Kang et al., 2024; Hooper et al., 2024; He et al., 2024; Su et al., 2025). Another key line incorporates **outlier handling**, by storing in higher bitwidths (Dong et al., 2024b; Kang et al., 2024; Hooper et al., 2024; Zhao et al., 2024b; Zandieh et al., 2025; Su et al., 2025), or by applying dedicated techniques shown in Tab. 2.

Recent advances explore vector quantization (VQ), which compresses groups using codebooks

Table 2: Summary of representative KV cache quantization (q.) methods. The "avg. bits" column shows the average
bitwidth per KV element, based on the authors' recommended or main experimental configuration(s).

Method	Gra Keys	nularity Values	Prec. mode	Important region	Outlier handling	Avg. bits
C 4.0 (W + 1 2022)				region		
SmoothQuant (Xiao et al., 2023)		nnel-wise	Fixed	_	Smoothing via scaling	8
FlexGen (Sheng et al., 2023)		up-wise	Fixed	_	_	4
WKVQuant (Yue et al., 2024)	2D (char	nnel & token)	Mixed	Current token	Dynamic token-wise q.	
MiKV (Yang et al., 2024b)	Tok	en-wise	Mixed	Existing policy	Outlier balancing	~4
QAQ (Dong et al., 2024b)	Tok	en-wise	Mixed	Attention-aware	Sparse matrix (FP16)	1.8 - 2.7
Atom (Zhao et al., 2024b)	Gro	up-wise	Mixed	Outlier channels	Selective high-bits	4.25
KIVI (Liu et al., 2024f)	Channel-wise	Token-wise	Mixed	Recent tokens	Channel-wise confining	~2
CacheGen (Liu et al., 2024e)	Layer-wise		Mixed	Shallow layers	_	1.9-2.9
DecoQuant (Liu et al., 2024d)	Decomposed-tensor-wise		Mixed	Small tensors	Tensor decomposition	4
GEAR (Kang et al., 2024)	Channel-wise	Token-wise	Fixed	_	Sparse matrix (FP16)	4.4/5.0
SKVQ (Duanmu et al., 2024)	Gro	up-wise	Mixed	Init. & recent tokens	Clipped dynamic q.	2.25
KVQuant (Hooper et al., 2024)	Channel-wise	Token-wise	Mixed	First token	Sparse matrix (FP16)	4.3
CQ (Zhang et al., 2024b)	Token-wise	channel-group	Fixed	_	_	1.3
ZipCache (He et al., 2024)	Channel-wise C	hansep. token-wise	Mixed	Norm. attention	Channel-wise norm.	3.2
QJL (Zandieh et al., 2025)	Tok	en-wise	Fixed	_	Selective high-bits	3/5
VQ-LLM (Liu et al., 2025b)	Group-wise	e (configurable)	Fixed	_		2/4
SQuat (Wang et al., 2025)	Block-wise	Token-wise	Fixed	_	_	3.1
QoQ (Lin et al., 2025b)	Char	nnel-wise	Fixed	-	Smooth attention	4
NSNQuant (Son et al., 2025)	Token-	wise vector	Fixed	_	Token-wise norm.	1.2/2.2
CommVQ (Li et al., 2025)	Token-	wise vector	Fixed	-	-	2
OTT (Su et al., 2025)	Channel-wise	Token-wise	Mixed	Outlier & recent tokens	Full precision	2.5

to capture inter-element correlations. CQ (Zhang et al., 2024b) couples channels and learns centroids for 1-bit quantization. VQ-LLM (Liu et al., 2025b) and CommVQ (Li et al., 2025) cut decode overhead via fused VQ kernels and RoPE-commutative codebooks. NSNQuant (Son et al., 2025) aligns KV distributions to avoid calibration errors.

5.1.2 KV Cache Low-rank Approximation

Low-rank compression exploits hidden-dimension redundancy by factorizing tensors into compact components. One family directly **projects KV tensors** to lower dimensions. FDC (Zhang and Shen, 2025) and xKV (Chang et al., 2025a) use singular value decomposition (SVD) on full KV caches, while ShadowKV (Sun et al., 2024) focuses on keys only. MatryoshkaKV (Lin et al., 2025a) learns orthogonal projections optimized via knowledge distillation. Similarly, LESS (Dong et al., 2024a) learns a low-rank residual to recover omitted signal.

To reduce projection cost, others **factorize KV projection weights**. LoRC (Zhang et al., 2024a) uses a progressive strategy guided by layer sensitivity. HeadKV (Yu et al., 2024b) and Palu (Chang et al., 2025b) group KV heads for SVD. Re-CalKV (Yan et al., 2025) decouples keys and values: it applies head-wise reordering for keys and offline calibration and matrix fusion for values.

5.1.3 KV Cache Structural Compression

Unlike value-level compression, structural compression modifies the organization of KV entries. **Structural sparsification** prunes redundant KV

structures. CHAI (Agarwal et al., 2024) clusters similar attention heads and computes attention only for representatives. ThinK (Xu et al., 2025b) prunes channels in keys via a query-driven criterion. ClusterAttn (Zhang et al., 2025b) clusters intrinsic attention patterns in the prompt after prefill.

Structural merging fuses KV entries into shared forms. Intra-layer methods merge similar tokens within a layer using attention scores (Zhang et al., 2024c), cosine similarity (Wan et al., 2024), or clustering with Gaussian-weighted fusion (Wang et al., 2024b). For cross-layer redundancy, Sub-Gen (Zandieh et al., 2024) and MiniCache (Liu et al., 2024a) merge entries from adjacent similar layers, while KVSharer (Yang et al., 2024c) surprisingly finds dissimilar-layer sharing more effective.

5.2 KV Cache Retention Management

Going beyond representations, this section focuses on mechanisms that manage the retention of the KV cache (*KVRM*) during serving, such as efficient allocation, reuse, and eviction strategies.

5.2.1 KV Cache Allocation and Reuse

Recent efforts explore dynamic KV allocation and reuse to improve memory efficiency. **Structure-aware** methods redesign KV cache layouts for flexible allocation and reuse. Inspired by OS virtual memory, PagedAttention (Kwon et al., 2023) uses fixed-size pages with logical-to-physical mapping, enabling lazy allocation and prefix reuse. Extending this design, vTensor (Xu et al., 2024a) decouples computation from defragmentation, and vAt-

tention (Prabhu et al., 2025) decouples virtual and physical memory allocation. ChunkAttention (Ye et al., 2024) and RadixAttention (Zheng et al., 2024) structure KV caches as prefix and radix trees for system prompt reuse. PromptCache (Gim et al., 2024) proposes the modular reuse of KV caches beyond the prefix. FastSwitch (Shen et al., 2024) adopts block group-based allocation and enables multi-turn reuse. LeanKV (Zhang et al., 2025c) performs parallel KV compaction to recycle fragmented memory. FlashInfer (Ye et al., 2025) presents a block-sparse and composable format.

Semantics-guided methods utilize semantic signals to guide KV allocation and reuse. LazyLLM (Fu et al., 2024) selectively materializes KV only for critical tokens, caching others for possible recomputation. MemServe (Hu et al., 2024a) extends this idea to distributed serving with a MemPool system, enabling reuse across and within requests.

5.2.2 KV Cache Eviction

KV cache eviction discards less critical KV entries (i.e., tokens) based on certain rules. **Static methods**, applied during or after prefill, keep the KV cache fixed during decoding (Ge et al., 2024; Li et al., 2024c; Zhou et al., 2024a). For example, SnapKV (Li et al., 2024c) identifies informative tokens via an observation window. Yet, such methods cannot adapt to importance shifts during decoding.

Dynamic eviction methods assess importance and perform eviction during decoding, often based on attention scores (Zhang et al., 2023; Liu et al., 2023; Oren et al., 2024) or heuristics (Adnan et al., 2024; Xiao et al., 2024b; Zhao et al., 2024a; Chen et al., 2025a). Beyond such metrics, RoCo (Ren and Zhu, 2024) and NACL (Chen et al., 2024) explore variance- or random-based strategies to mitigate bias in local attention statistics. VATP (Guo et al., 2024) incorporates value L_1 -norms to complement attention scores. Devoto et al. (2024) similarly correlated L_2 -norm with attention scores.

Observing layer-wise differences in attention, layer-aware methods assign retention budgets accordingly. PyramidKV (Cai et al., 2024) and PyramidInfer (Yang et al., 2024a) favor shallow layers, while D₂O (Wan et al., 2024) and CAKE (Qin et al., 2025) allocate by attention density or layer preferences. LaCache (Shi et al., 2025) implicitly allocates layer budgets through a ladder-shaped retention pattern and iterative eviction. At finer granularity, Ada-KV (Feng et al., 2024) reallocates budgets across heads within each layer.

6 Observations and Open Challenges

Here, we identify observations from two complementary lenses: (i) a behavior \times objective matrix and (ii) a behavior-behavior synergy network, which naturally motivate our open challenges.

Table 3 (behavior \times objective matrix) marks each behavior's impact on serving objectives as direct (\odot) or indirect (\bigcirc); stars (\star) on direct cells statistically flag $\geq 70\%$ of papers reporting such gains. Side bars show research density. Objectives cover latency, throughput, GPU memory, interconnect I/O, and energy. We also include quality impact \downarrow to capture degradation as a trade-off. Figure 6 (behavior-behavior synergy network) visualizes cross-behavior co-occurrence, with edge thickness proportional to normalized weights. Lowscore edges are omitted, and computation details are in App. E. Key observations are as follows.

- O1. Structural behaviors receive the most attention and dominate GPU memory savings, while other behaviors yield savings indirectly (e.g., via migration or reuse). This indicates a community default to prioritize memory efficiency.
- **O2.** Temporal behaviors act most directly on latency and throughput, since KVS, OVLP, and HAE map cleanly to reductions in scheduling stalls, pipeline bubbles, and device under-utilization. However, tail latency reporting is sparse.
- **O3.** Spatial methods primarily target interconnect I/O, often paired with OVLP. Their core focus is KV transfer, and by overlapping it with compute, they effectively hide transfer latency.
- **O4. Energy is under-explored,** although many methods reduce memory or compute intensity that should translate to energy benefits.
- **O5. Quality loss is universal.** Temporal methods risk inconsistent request handling; spatial methods risk missed KV data; and structural methods directly reduce KV precision. The practical question is to ensure such degradation is controllable and aligned with service-level objectives (SLOs).
- **O6. HAE–CDO** is the strongest cross-dimension synergy. Compute layouts that exploit device heterogeneity often co-design with KV colocating or transfer, yielding joint gains in utilization and I/O. **O7. KVCC** remains isolated, despite its popularity. Structural synergies are sparse overall, suggesting a missed opportunity for co-design.

The above observations reveal both progress and gaps of current sKis research, which motivate the next set of system-level open challenges.

Table 3: Behavior \times objective matrix of sKis methods. Side bars encode research density (rows/columns). Cells mark relevance levels (\bullet = direct, \bigcirc = indirect) and high-prevalence flags (\star : $\geq 70\%$ of papers report gains).

Behaviors	Mean latency	Tail latency	Through- put	GPU memory	Inter- connect I/O	Energy /power	Quality impact ↓	Row density
KV-centric scheduling	•*			0	0	0	0	
Pipelining and overlapping	•*	0	•*		\circ	0	0	
Hardware-aware execution	● *		● ★				0	
Memory hierarchy KV orchestration	0	0	\circ	0	•*			
Compute device KV orchestration	0	\circ	● ★	\circ				
KV cache compression	0	\circ	\circ	•*				
KV cache retention management	0	0	0	•*	•			
Column density								
Behavior dimension: temporal, spatial, structural.	lensity	0-10 11-2	20 21–30	31+	lumn density	0–20 21–	40 41–60	61+

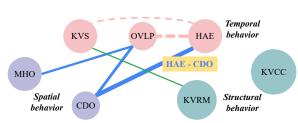


Figure 6: Behavior-behavior synergy network. Node size reflects research density; edge thickness scales with behavior co-occurrence. Low-score edges are omitted.

C1. SLO-driven tail control ← O2. SLOs are critical to LLM serving, with tail latency dominating user experience (Dean and Barroso, 2013; Wang et al., 2024c), yet most systems omit detailed reporting. Under long contexts and bursts, KV generation, placement, migration, and compression may interfere and trigger SLO violations. The challenge is to attribute SLO violations to concrete KV behaviors and paths, motivating studies on standardized preemption and degradation semantics, e.g., safe interrupt points and bounded decompression costs, to make tail outcomes controllable.

C2. Energy-aware sKis \leftarrow O4. With surging data center demand, sKis should be energy-aware, but energy is rarely reported or optimized. Future research could integrate power profiling into runtime decisions, establish serving-time energy models, and jointly optimize energy-latency-quality under power constraints. Another possible direction is to study energy-friendly KV granularities and layouts. C3. Trustworthy and efficient sKis \leftarrow O5. The success of LLM serving must ensure not only quality but also trustworthiness (Han et al., 2025). Efficiency optimizations often consider quality loss, yet trustworthiness is rarely measured or attributed. The challenge lies in identifying behaviors that degrade trust. For example, temporal asynchrony may expose stale KV and add nondeterminism, harming reliability; cross-tier migration risks KV residue

or plaintext transfer, harming privacy; and KVCC carries inherent risk of discarding rare but critical signals, harming robustness.

C4. Generalizable HAE-CDO ← O6. While HAE and CDO form the strongest synergy, policies are often tailored to specific fabric or single-tenant settings. Future directions include making such synergy portable across heterogeneous topologies (e.g., NVLink, NVSwitch, PCIe, CXL) and adaptive to multi-tenant settings with mixed workloads.

C5. Co-optimization and intermediate semantics ← O7. Most sKis optimize behaviors in isolation, despite their interactions under bandwidth and latency pressure. Future studies could explore share-budget cost models for joint optimization, For instance, to co-decide eviction, offload, and prefetch given predicted reuse, success probability, and I/O contention. Another useful direction is to introduce fine-grained intermediate semantics for behaviors, e.g., degrade-then-offload with graded surrogates and soft-evict with per-token grace.

C6. Unified benchmarks for sKis. We summarize in App. F a review of existing LLM inference benchmarking tools and practices. Despite this review, we find inconsistent metric definitions and measurements across tools. It is worth developing standardized benchmarks that unify performance metrics and service-level indicators with clear definitions and window semantics.

7 Conclusion

This survey presents a systematic overview of sKis, offering a system behavior-oriented taxonomy covering temporal, spatial, and structural dimensions. By cross-analyzing behavior-objective impacts and behavior-behavior synergies, we reveal overlooked regions, potential synergies, and open challenges. We hope this survey inspires continued exploration toward efficient and trustworthy LLM serving.

Limitations

This paper offers a comprehensive review and summary of current methods in the area of systemaware KV cache optimization. However, given the extensive body of related work and the rapidly evolving nature of this research area, we may have overlooked some equally valuable contributions. We tried to include all relevant studies and references wherever feasible. Additionally, this survey conducts no new experiments. Our claims synthesize results reported in public papers and open-source implementations, primarily under mainstream platforms and common configurations, which may constrain the generality of our conclusions. Finally, we outline several KV-centric research directions to improve the efficiency in LLM serving, including SLO-first tail-latency control, energy-aware sKis, trustworthy sKis, generalizable HAE-CDO, co-optimization and intermediate semantics, and unified benchmarks. We plan to leave these aspects for future work.

References

- Muhammad Adnan, Akhil Arunkumar, Gaurav Jain, Prashant J Nair, Ilya Soloveychik, and Purushotham Kamath. 2024. Keyformer: KV cache reduction through key tokens selection for efficient generative inference. *Proceedings of Machine Learning and Systems*, 6:114–127.
- Saurabh Agarwal, Bilge Acun, Basil Hosmer, Mostafa Elhoushi, Yejin Lee, Shivaram Venkataraman, Dimitris Papailiopoulos, and Carole-Jean Wu. 2024. CHAI: Clustered head attention for efficient LLM inference. In *International Conference on Machine Learning*, pages 291–312. PMLR.
- Reza Yazdani Aminabadi, Samyam Rajbhandari, Ammar Ahmad Awan, Cheng Li, Du Li, Elton Zheng, Olatunji Ruwase, Shaden Smith, Minjia Zhang, Jeff Rasley, and Yuxiong He. 2022. DeepSpeed-Inference: enabling efficient inference of transformer models at unprecedented scale. In *International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15. IEEE.
- Artem Babenko and Victor Lempitsky. 2014. Additive quantization for extreme vector compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 931–938.
- Guangji Bai, Zheng Chai, Chen Ling, Shiyu Wang, Jiaying Lu, Nan Zhang, Tingwei Shi, Ziyang Yu, Mengdan Zhu, Yifei Zhang, Carl Yang, Yue Cheng, and Liang Zhao. 2024. Beyond efficiency: A systematic survey of resource-efficient large language models. arXiv preprint arXiv:2401.00625.

- 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, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong, Baobao Chang, Junjie Hu, and Wen Xiao. 2024. PyramidKV: Dynamic KV cache compression based on pyramidal information funneling. *arXiv preprint arXiv:2406.02069*.
- Chi-Chih Chang, Chien-Yu Lin, Yash Akhauri, Wei-Cheng Lin, Kai-Chiang Wu, Luis Ceze, and Mohamed S Abdelfattah. 2025a. xKV: Cross-layer SVD for KV-cache compression. *arXiv* preprint *arXiv*:2503.18893.
- Chi-Chih Chang, Wei-Cheng Lin, Chien-Yu Lin, Chong-Yan Chen, Yu-Fang Hu, Pei-Shuo Wang, Ning-Chi Huang, Luis Ceze, Mohamed S Abdelfattah, and Kai-Chiang Wu. 2025b. Palu: KV-cache compression with low-rank projection. In *International Conference on Learning Representations*.
- Guoxuan Chen, Han Shi, Jiawei Li, Yihang Gao, Xiaozhe Ren, Yimeng Chen, Xin Jiang, Zhenguo Li, Weiyang Liu, and Chao Huang. 2025a. SepLLM: Accelerate large language models by compressing one segment into one separator. In *International Conference on Machine Learning*.
- Yilong Chen, Guoxia Wang, Junyuan Shang, Shiyao Cui, Zhenyu Zhang, Tingwen Liu, Shuohuan Wang, Yu Sun, Dianhai Yu, and Hua Wu. 2024. NACL: A general and effective KV cache eviction framework for LLM at inference time. In *Annual Meeting of the Association for Computational Linguistics*, pages 7913–7926.
- Zhuoming Chen, Ranajoy Sadhukhan, Zihao Ye, Yang Zhou, Jianyu Zhang, Niklas Nolte, Yuandong Tian, Matthijs Douze, Léon Bottou, Zhihao Jia, and Beidi Chen. 2025b. MagicPIG: LSH sampling for efficient LLM generation. In *International Conference on Learning Representations*.
- Krishna Teja Chitty-Venkata, Siddhisanket Raskar, Bharat Kale, Farah Ferdaus, Aditya Tanikanti, Ken Raffenetti, Valerie Taylor, Murali Emani, and Venkatram Vishwanath. 2024. LLM-Inference-Bench: Inference benchmarking of large language models on AI accelerators. In Workshops of the International Conference for High Performance Computing, Networking, Storage and Analysis, pages 1362–1379. IEEE.
- Jeffrey Dean and Luiz André Barroso. 2013. The tail at scale. *Communications of the ACM*, 56(2):74–80.

- Alessio Devoto, Yu Zhao, Simone Scardapane, and Pasquale Minervini. 2024. A simple and effective L_2 norm-based strategy for KV cache compression. In *Conference on Empirical Methods in Natural Language Processing*, pages 18476–18499.
- Yiran Ding, Li Lyna Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang, and Mao Yang. 2024. LongRoPE: Extending LLM context window beyond 2 million tokens. *arXiv* preprint arXiv:2402.13753.
- Harry Dong, Xinyu Yang, Zhenyu Zhang, Zhangyang Wang, Yuejie Chi, and Beidi Chen. 2024a. Get more with LESS: Synthesizing recurrence with KV cache compression for efficient LLM inference. *arXiv* preprint arXiv:2402.09398.
- Shichen Dong, Wen Cheng, Jiayu Qin, and Wei Wang. 2024b. QAQ: Quality adaptive quantization for LLM KV cache. *arXiv preprint arXiv:2403.04643*.
- Yanhao Dong, Yubo Miao, Weinan Li, Xiao Zheng, Chao Wang, and Feng Lyu. 2025. Accelerating LLM inference throughput via asynchronous KV cache prefetching. *arXiv preprint arXiv:2504.06319*.
- Jiangfei Duan, Runyu Lu, Haojie Duanmu, Xiuhong Li, Xingcheng Zhang, Dahua Lin, Ion Stoica, and Hao Zhang. 2024. MuxServe: Flexible spatial-temporal multiplexing for multiple LLM serving. In *International Conference on Machine Learning*, pages 11905–11917. PMLR.
- Haojie Duanmu, Zhihang Yuan, Xiuhong Li, Jiangfei Duan, Xingcheng Zhang, and Dahua Lin. 2024. SKVQ: Sliding-window key and value cache quantization for large language models. *arXiv preprint arXiv*:2405.06219.
- Yuan Feng, Junlin Lv, Yukun Cao, Xike Xie, and S Kevin Zhou. 2024. Ada-KV: Optimizing KV cache eviction by adaptive budget allocation for efficient LLM inference. *arXiv* preprint arXiv:2407.11550.
- Qichen Fu, Minsik Cho, Thomas Merth, Sachin Mehta, Mohammad Rastegari, and Mahyar Najibi. 2024. LazyLLM: Dynamic token pruning for efficient long context LLM inference. *arXiv preprint arXiv:2407.14057*.
- Bin Gao, Zhuomin He, Puru Sharma, Qingxuan Kang, Djordje Jevdjic, Junbo Deng, Xingkun Yang, Zhou Yu, and Pengfei Zuo. 2024. Cost-efficient large language model serving for multi-turn conversations with Cached Attention. In *USENIX Annual Technical Conference*, pages 111–126.
- Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. 2024. Model tells you what to discard: Adaptive KV cache compression for LLMs. In *International Conference on Learning Representations*.

- In Gim, Guojun Chen, Seung-seob Lee, Nikhil Sarda, Anurag Khandelwal, and Lin Zhong. 2024. Prompt Cache: Modular attention reuse for low-latency inference. *Proceedings of Machine Learning and Systems*, 6:325–338.
- Robert Gray. 1984. Vector quantization. *IEEE Assp Magazine*, 1(2):4–29.
- Xiangming Gu, Tianyu Pang, Chao Du, Qian Liu, Fengzhuo Zhang, Cunxiao Du, Ye Wang, and Min Lin. 2025. When attention sink emerges in language models: An empirical view. In *International Conference on Learning Representations*.
- Zhiyu Guo, Hidetaka Kamigaito, and Taro Watanabe. 2024. Attention score is not all you need for token importance indicator in KV cache reduction: Value also matters. In *Conference on Empirical Methods in Natural Language Processing*, pages 21158–21166.
- Bo Han, Jiangchao Yao, Tongliang Liu, Bo Li, Sanmi Koyejo, and Feng Liu. 2025. Trustworthy machine learning: From data to models. *Foundations and Trends*® *in Privacy and Security*, 7(2-3):74–246.
- Jiaao He and Jidong Zhai. 2024. FastDecode: Highthroughput GPU-efficient LLM serving using heterogeneous pipelines. arXiv preprint arXiv:2403.11421.
- Yefei He, Luoming Zhang, Weijia Wu, Jing Liu, Hong Zhou, and Bohan Zhuang. 2024. ZipCache: Accurate and efficient KV cache quantization with salient token identification. In *Advances in Neural Information Processing Systems*, volume 37, pages 68287–68307.
- Yintao He, Haiyu Mao, Christina Giannoula, Mohammad Sadrosadati, Juan Gómez-Luna, Huawei Li, Xiaowei Li, Ying Wang, and Onur Mutlu. 2025. PAPI: Exploiting dynamic parallelism in large language model decoding with a processing-in-memory-enabled computing system. In ACM International Conference on Architectural Support for Programming Languages and Operating Systems, pages 766–782.
- Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Yakun S Shao, Kurt Keutzer, and Amir Gholami. 2024. KVQuant: Towards 10 million context length LLM inference with KV cache quantization. *Advances in Neural Information Processing Systems*, 37:1270–1303.
- Cunchen Hu, HeYang Huang, Junhao Hu, Jiang Xu, Xusheng Chen, Tao Xie, Chenxi Wang, Sa Wang, Yungang Bao, Ninghui Sun, and Yizhou Shan. 2024a. MemServe: Context caching for disaggregated LLM serving with elastic memory pool. *arXiv preprint arXiv:2406.17565*.
- Cunchen Hu, HeYang Huang, Liangliang Xu, Xusheng Chen, Jiang Xu, Shuang Chen, Hao Feng, Chenxi Wang, Sa Wang, Yungang Bao, Ninghui Sun, and Yizhou Shan. 2024b. Inference without interference: Disaggregate LLM inference for mixed downstream workloads. *arXiv preprint arXiv:2401.11181*.

- Herve Jegou, Matthijs Douze, and Cordelia Schmid. 2010. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence*, 33(1):117–128.
- Chaoyi Jiang, Lei Gao, Hossein Entezari Zarch, and Murali Annavaram. 2025a. KVPR: Efficient LLM inference with i/o-aware KV cache partial recomputation. In *Annual Meeting of the Association for Computational Linguistics*.
- Jiantong Jiang, Zeyi Wen, Atif Mansoor, and Ajmal Mian. 2024. Fast inference for probabilistic graphical models. In USENIX Annual Technical Conference, pages 95–110.
- Jiantong Jiang, Zeyi Wen, and Ajmal Mian. 2022. Fast parallel Bayesian network structure learning. In *IEEE International Parallel and Distributed Processing Symposium*, pages 617–627. IEEE.
- Xuanlin Jiang, Yang Zhou, Shiyi Cao, Ion Stoica, and Minlan Yu. 2025b. Neo: Saving GPU memory crisis with CPU offloading for online LLM inference. *Proceedings of Machine Learning and Systems*.
- Chao Jin, Zili Zhang, Xuanlin Jiang, Fangyue Liu, Xin Liu, Xuanzhe Liu, and Xin Jin. 2024. RAGCache: Efficient knowledge caching for retrieval-augmented generation. *arXiv* preprint arXiv:2404.12457.
- Lena Jurkschat, Preetam Gattogi, Sahar Vahdati, and Jens Lehmann. 2025. BALI-a benchmark for accelerated language model inference. *IEEE Access*.
- Hao Kang, Qingru Zhang, Souvik Kundu, Geonhwa Jeong, Zaoxing Liu, Tushar Krishna, and Tuo Zhao. 2024. GEAR: An efficient KV cache compression recipe for near-lossless generative inference of LLM. arXiv preprint arXiv:2403.05527.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with PagedAttention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626.
- Wonbeom Lee, Jungi Lee, Junghwan Seo, and Jaewoong Sim. 2024. InfiniGen: Efficient generative inference of large language models with dynamic KV cache management. In *USENIX Symposium on Operating Systems Design and Implementation*, pages 155–172.
- Baolin Li, Yankai Jiang, Vijay Gadepally, and Devesh Tiwari. 2024a. LLM inference serving: Survey of recent advances and opportunities. *arXiv preprint arXiv:2407.12391*.
- Haoyang Li, Yiming Li, Anxin Tian, Tianhao Tang, Zhanchao Xu, Xuejia Chen, Nicole Hu, Wei Dong, Qing Li, and Lei Chen. 2024b. A survey on large language model acceleration based on KV cache management. *arXiv preprint arXiv:2412.19442*.

- Junyan Li, Yang Zhang, Muhammad Yusuf Hassan, Talha Chafekar, Tianle Cai, Zhile Ren, Pengsheng Guo, Foroozan Karimzadeh, Chong Wang, and Chuang Gan. 2025. CommVQ: Commutative vector quantization for KV cache compression. In *Interna*tional Conference on Machine Learning.
- Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. 2024c. SnapKV: LLM knows what you are looking for before generation. *Advances in Neural Information Processing Systems*, 37:22947–22970.
- Bin Lin, Tao Peng, Chen Zhang, Minmin Sun, Lanbo Li, Hanyu Zhao, Wencong Xiao, Qi Xu, Xiafei Qiu, Shen Li, Zhigang Ji, Yong Li, and Wei Lin. 2024. Infinite-LLM: Efficient LLM service for long context with distattention and distributed KVCache. *arXiv* preprint arXiv:2401.02669.
- Bokai Lin, Zihao Zeng, Zipeng Xiao, Siqi Kou, Tianqi Hou, Xiaofeng Gao, Hao Zhang, and Zhijie Deng. 2025a. MatryoshkaKV: Adaptive kv compression via trainable orthogonal projection. In *International Conference on Learning Representations*.
- Yujun Lin, Haotian Tang, Shang Yang, Zhekai Zhang, Guangxuan Xiao, Chuang Gan, and Song Han. 2025b. QServe: W4A8KV4 quantization and system codesign for efficient LLM serving. *Proceedings of Machine Learning and Systems*.
- Akide Liu, Jing Liu, Zizheng Pan, Yefei He, Gholamreza Haffari, and Bohan Zhuang. 2024a. MiniCache: KV cache compression in depth dimension for large language models. In *Advances in Neural Informa*tion Processing Systems, volume 37, pages 139997– 140031.
- Di Liu, Meng Chen, Baotong Lu, Huiqiang Jiang, Zhenhua Han, Qianxi Zhang, Qi Chen, Chengruidong Zhang, Bailu Ding, Kai Zhang, Chen Chen, Fan Yang, Yuqing Yang, and Lili Qiu. 2024b. Retrieval Attention: Accelerating long-context LLM inference via vector retrieval. arXiv preprint arXiv:2409.10516.
- Guangda Liu, Chengwei Li, Jieru Zhao, Chenqi Zhang, and Minyi Guo. 2024c. ClusterKV: Manipulating LLM KV cache in semantic space for recallable compression. arXiv preprint arXiv:2412.03213.
- Peiyu Liu, Ze-Feng Gao, Xin Zhao, Yipeng Ma, Tao Wang, and Ji-Rong Wen. 2024d. Unlocking data-free low-bit quantization with matrix decomposition for KV cache compression. In *Annual Meeting of the Association for Computational Linguistics*, pages 2430–2440.
- Yanyu Liu, Jingying Fu, Sixiang Liu, Yitian Zou, You Fu, Jiehan Zhou, and Shouhua Zhang. 2025a. KV cache compression for inference efficiency in LLMs: A review. *arXiv preprint arXiv:2508.06297*.
- Yuhan Liu, Hanchen Li, Yihua Cheng, Siddhant Ray, Yuyang Huang, Qizheng Zhang, Kuntai Du, Jiayi

- Yao, Shan Lu, Ganesh Ananthanarayanan, Michael Maire, Henry Hoffmann, Ari Holtzman, and Junchen Jiang. 2024e. CacheGen: KV cache compression and streaming for fast large language model serving. In *Proceedings of the ACM SIGCOMM 2024 Conference*, pages 38–56.
- Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios Kyrillidis, and Anshumali Shrivastava. 2023. Scissorhands: Exploiting the persistence of importance hypothesis for LLM KV cache compression at test time. Advances in Neural Information Processing Systems, 36:52342–52364.
- Zihan Liu, Xinhao Luo, Junxian Guo, Wentao Ni, Yangjie Zhou, Yue Guan, Cong Guo, Weihao Cui, Yu Feng, Minyi Guo, Yuhao Zhu, Minjia Zhang, Jingwen Leng, and Chen Jin. 2025b. VQ-LLM: Highperformance code generation for vector quantization augmented LLM inference. In *IEEE International Symposium on High Performance Computer Architecture*, pages 1496–1509. IEEE.
- Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. 2024f. KIVI: A tuning-free asymmetric 2bit quantization for KV cache. In *International Conference on Machine Learning*, pages 32332–32344. PMLR.
- Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Hongyi Jin, Tianqi Chen, and Zhihao Jia. 2023. Towards efficient generative large language model serving: A survey from algorithms to systems. arXiv preprint arXiv:2312.15234.
- NVIDIA. 2023. TensorRT-LLM. https://github.com/NVIDIA/TensorRT-LLM.
- NVIDIA. 2025a. GenAI-Perf. NVIDIA Docs.
- NVIDIA. 2025b. NVIDIA NIM LLMs benchmarking. *NVIDIA Docs*.
- OpenAI. 2023. GPT-4 technical report. arXiv preprint arXiv:2303.08774.
- Matanel Oren, Michael Hassid, Nir Yarden, Yossi Adi, and Roy Schwartz. 2024. Transformers are multistate RNNs. In *Conference on Empirical Methods in Natural Language Processing*, pages 18724–18741.
- Xiurui Pan, Endian Li, Qiao Li, Shengwen Liang, Yizhou Shan, Ke Zhou, Yingwei Luo, Xiaolin Wang, and Jie Zhang. 2024. InstInfer: In-storage attention offloading for cost-effective long-context LLM inference. arXiv preprint arXiv:2409.04992.
- Daon Park and Bernhard Egger. 2024. Improving throughput-oriented LLM inference with CPU computations. In *International Conference on Parallel Architectures and Compilation Techniques*, pages 233–245.

- Daon Park, Sungbin Jo, and Bernhard Egger. 2023. Improving throughput-oriented generative inference with CPUs. In *ACM SIGOPS Asia-Pacific Workshop on Systems*, pages 37–42.
- Jaehyun Park, Jaewan Choi, Kwanhee Kyung, Michael Jaemin Kim, Yongsuk Kwon, Nam Sung Kim, and Jung Ho Ahn. 2024. AttAcc! unleashing the power of PIM for batched transformer-based generative model inference. In ACM International Conference on Architectural Support for Programming Languages and Operating Systems, pages 103–119.
- Pratyush Patel, Esha Choukse, Chaojie Zhang, Aashaka Shah, Íñigo Goiri, Saeed Maleki, and Ricardo Bianchini. 2024. Splitwise: Efficient generative LLM inference using phase splitting. In *ACM/IEEE Annual International Symposium on Computer Architecture*, pages 118–132. IEEE.
- Sundar Pichai and Demis Hassabis. 2024. Our next-generation model: Gemini 1.5. *Google Blog*.
- Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. 2023. Efficiently scaling transformer inference. *Proceedings of Machine Learning and Systems*, 5:606–624.
- Ramya Prabhu, Ajay Nayak, Jayashree Mohan, Ramachandran Ramjee, and Ashish Panwar. 2025. vAttention: Dynamic memory management for serving LLMs without PagedAttention. In ACM International Conference on Architectural Support for Programming Languages and Operating Systems, pages 1133–1150.
- Ruoyu Qin, Zheming Li, Weiran He, Mingxing Zhang, Yongwei Wu, Weimin Zheng, and Xinran Xu. 2024. Mooncake: A KVCache-centric disaggregated architecture for LLM serving. *arXiv* preprint *arXiv*:2407.00079.
- Ziran Qin, Yuchen Cao, Mingbao Lin, Wen Hu, Shixuan Fan, Ke Cheng, Weiyao Lin, and Jianguo Li. 2025. CAKE: Cascading and adaptive KV cache eviction with layer preferences. In *International Conference on Learning Representations*.
- Haoran Qiu, Weichao Mao, Archit Patke, Shengkun Cui, Saurabh Jha, Chen Wang, Hubertus Franke, Zbigniew Kalbarczyk, Tamer Başar, and Ravishankar K Iyer. 2024. Power-aware deep learning model serving with μ-serve. In *USENIX Annual Technical Conference*, pages 75–93.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2020. ZeRO: Memory optimizations toward training trillion parameter models. In *International Conference for High Performance Computing*, *Networking, Storage and Analysis*, pages 1–16. IEEE.
- Ray. 2024. LLMPerf. https://github.com/ray-project/llmperf.
- Vijay Janapa Reddi, Christine Cheng, David Kanter, Peter Mattson, Guenther Schmuelling, Carole-Jean Wu, Brian Anderson, Maximilien Breughe, Mark Charlebois, William Chou, and 1 others. 2020. MLPerf inference benchmark. In *ACM/IEEE Annual International Symposium on Computer Architecture*, pages 446–459. IEEE.
- Siyu Ren and Kenny Q Zhu. 2024. On the efficacy of eviction policy for key-value constrained generative language model inference. *arXiv* preprint *arXiv*:2402.06262.
- Luka Ribar, Ivan Chelombiev, Luke Hudlass-Galley, Charlie Blake, Carlo Luschi, and Douglas Orr. 2024. SparQ Attention: Bandwidth-efficient LLM inference. In *International Conference on Machine Learning*, pages 42558–42583. PMLR.
- Rana Shahout, Cong Liang, Shiji Xin, Qianru Lao, Yong Cui, Minlan Yu, and Michael Mitzenmacher. 2024. Fast inference for augmented large language models. *arXiv preprint arXiv:2410.18248*.
- Ao Shen, Zhiyao Li, and Mingyu Gao. 2024. Fastswitch: Optimizing context switching efficiency in fairness-aware large language model serving. arXiv preprint arXiv:2411.18424.
- Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, Beidi Chen, Percy Liang, Christopher Ré, Ion Stoica, and Ce Zhang. 2023. FlexGen: High-throughput generative inference of large language models with a single gpu. In *International Conference on Machine Learning*, pages 31094–31116. PMLR.
- Dachuan Shi, Yonggan Fu, Xiangchi Yuan, Zhongzhi Yu, Haoran You, Sixu Li, Xin Dong, Jan Kautz, Pavlo Molchanov, and Yingyan Lin. 2025. LaCache: Ladder-shaped KV caching for efficient long-context modeling of large language models. In *International Conference on Machine Learning*.
- Luohe Shi, Hongyi Zhang, Yao Yao, Zuchao Li, and Hai Zhao. 2024. Keep the cost down: A review on methods to optimize LLM's KV-cache consumption. *arXiv preprint arXiv:2407.18003*.
- Prajwal Singhania, Siddharth Singh, Shwai He, Soheil Feizi, and Abhinav Bhatele. 2024. Loki: Lowrank keys for efficient sparse attention. *Advances in Neural Information Processing Systems*, 37:16692–16723.

- Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. 2017. SmoothGrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*.
- Donghyun Son, Euntae Choi, and Sungjoo Yoo. 2025. NSNQuant: A double normalization approach for calibration-free low-bit vector quantization of KV cache. *arXiv* preprint arXiv:2505.18231.
- Vikranth Srivatsa, Zijian He, Reyna Abhyankar, Dongming Li, and Yiying Zhang. 2024. Preble: Efficient distributed prompt scheduling for llm serving. In *International Conference on Learning Representations*.
- Foteini Strati, Sara McAllister, Amar Phanishayee, Jakub Tarnawski, and Ana Klimovic. 2024. Déjàvu: KV-cache streaming for fast, fault-tolerant generative LLM serving. In *International Conference on Machine Learning*, pages 46745–46771.
- Yi Su, Yuechi Zhou, Quantong Qiu, Juntao Li, Qingrong Xia, Ping Li, Xinyu Duan, Zhefeng Wang, and Min Zhang. 2025. Accurate KV cache quantization with outlier tokens tracing. In *Annual Meeting of the Association for Computational Linguistics*.
- Hanshi Sun, Li-Wen Chang, Wenlei Bao, Size Zheng, Ningxin Zheng, Xin Liu, Harry Dong, Yuejie Chi, and Beidi Chen. 2024. ShadowKV: KV cache in shadows for high-throughput long-context LLM inference. arXiv preprint arXiv:2410.21465.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319—3328. PMLR.
- Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. 2024. Quest: queryaware sparsity for efficient long-context LLM inference. In *International Conference on Machine Learn*ing, pages 47901–47911.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aur'elien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. LLaMA: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cris tian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023b. LLaMA 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

- Zhongwei Wan, Xinjian Wu, Yu Zhang, Yi Xin, Chaofan Tao, Zhihong Zhu, Xin Wang, Siqi Luo, Jing Xiong, and Mi Zhang. 2024. D2O: Dynamic discriminative operations for efficient generative inference of large language models. *arXiv preprint arXiv:2406.13035*.
- Guanhua Wang, Zhuang Liu, Brandon Hsieh, Siyuan Zhuang, Joseph Gonzalez, Trevor Darrell, and Ion Stoica. 2021. sensAI: Convnets decomposition via class parallelism for fast inference on live data. *Proceedings of Machine Learning and Systems*, 3:664–679.
- Hao Wang, Ligong Han, Kai Xu, and Akash Srivastava. 2025. SQuat: Subspace-orthogonal KV cache quantization. *arXiv preprint arXiv:2503.24358*.
- Wenxiao Wang, Wei Chen, Yicong Luo, Yongliu Long, Zhengkai Lin, Liye Zhang, Binbin Lin, Deng Cai, and Xiaofei He. 2024a. Model compression and efficient inference for large language models: A survey. *arXiv preprint arXiv:2402.09748*.
- Zheng Wang, Boxiao Jin, Zhongzhi Yu, and Minjia Zhang. 2024b. Model tells you where to merge: Adaptive KV cache merging for LLMs on long-context tasks. *arXiv preprint arXiv:2407.08454*.
- Zhibin Wang, Shipeng Li, Yuhang Zhou, Xue Li, Rong Gu, Nguyen Cam-Tu, Chen Tian, and Sheng Zhong. 2024c. Revisiting SLO and goodput metrics in LLM serving. *arXiv preprint arXiv:2410.14257*.
- Bingyang Wu, Shengyu Liu, Yinmin Zhong, Peng Sun, Xuanzhe Liu, and Xin Jin. 2024a. LoongServe: Efficiently serving long-context large language models with elastic sequence parallelism. In *Proceedings of the ACM SIGOPS 30th Symposium on Operating Systems Principles*, pages 640–654.
- Bingyang Wu, Yinmin Zhong, Zili Zhang, Shengyu Liu, Fangyue Liu, Yuanhang Sun, Gang Huang, Xuanzhe Liu, and Xin Jin. 2023. Fast distributed inference serving for large language models. *arXiv preprint* arXiv:2305.05920.
- Wei Wu, Zhuoshi Pan, Chao Wang, Liyi Chen, Yunchu Bai, Tianfu Wang, Kun Fu, Zheng Wang, and Hui Xiong. 2024b. TokenSelect: Efficient long-context inference and length extrapolation for LLMs via dynamic token-level KV cache selection. *arXiv* preprint *arXiv*:2411.02886.
- Heming Xia, Zhe Yang, Qingxiu Dong, Peiyi Wang, Yongqi Li, Tao Ge, Tianyu Liu, Wenjie Li, and Zhifang Sui. 2024. Unlocking efficiency in large language model inference: A comprehensive survey of speculative decoding. In *Findings of the Association for Computational Linguistics ACL*, pages 7655–7671.
- Chaojun Xiao, Pengle Zhang, Xu Han, Guangxuan Xiao, Yankai Lin, Zhengyan Zhang, Zhiyuan Liu, and Maosong Sun. 2024a. InfLLM: Training-free long-context extrapolation for LLMs with an efficient context memory. *Advances in Neural Information Processing Systems*, 37:119638–119661.

- 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 *International Conference on Machine Learning*, pages 38087–38099. PMLR.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024b. Efficient streaming language models with attention sinks. In *International Conference on Learning Representations*.
- Wencong Xiao, Romil Bhardwaj, Ramachandran Ramjee, Muthian Sivathanu, Nipun Kwatra, Zhenhua Han, Pratyush Patel, Xuan Peng, Hanyu Zhao, Quanlu Zhang, Fan Yang, and Lidong Zhou. 2018. Gandiva: Introspective cluster scheduling for deep learning. In USENIX Symposium on Operating Systems Design and Implementation, pages 595–610.
- Yi Xiong, Hao Wu, Changxu Shao, Ziqing Wang, Rui Zhang, Yuhong Guo, Junping Zhao, Ke Zhang, and Zhenxuan Pan. 2024. LayerKV: Optimizing large language model serving with layer-wise KV cache management. *arXiv preprint arXiv:2410.00428*.
- Fangyuan Xu, Tanya Goyal, and Eunsol Choi. 2025a. RefreshKV: Updating small KV cache during longform generation. In Annual Meeting of the Association for Computational Linguistics, pages 24878– 24893.
- Jiale Xu, Rui Zhang, Cong Guo, Weiming Hu, Zihan Liu, Feiyang Wu, Yu Feng, Shixuan Sun, Changxu Shao, Yuhong Guo, Junping Zhao, Ke Zhang, Minyi Guo, and Jingwen Leng. 2024a. vTensor: Flexible virtual tensor management for efficient LLM serving. arXiv preprint arXiv:2407.15309.
- Mengwei Xu, Wangsong Yin, Dongqi Cai, Rongjie Yi, Daliang Xu, Qipeng Wang, Bingyang Wu, Yihao Zhao, Chen Yang, Shihe Wang, Qiyang Zhang, Zhenyan Lu, Li Zhang, Shangguang Wang, Yuanchun Li, Yunxin Liu, Xin Jin, and Xuanzhe Liu. 2024b. A survey of resource-efficient LLM and multimodal foundation models. *arXiv preprint arXiv:2401.08092*.
- Yuhui Xu, Zhanming Jie, Hanze Dong, Lei Wang, Xudong Lu, Aojun Zhou, Amrita Saha, Caiming Xiong, and Doyen Sahoo. 2025b. ThinK: Thinner key cache by query-driven pruning. In *International Conference on Learning Representations*.
- Xianglong Yan, Zhiteng Li, Tianao Zhang, Linghe Kong, Yulun Zhang, and Xiaokang Yang. 2025. Re-CalKV: Low-rank KV cache compression via head reordering and offline calibration. *arXiv preprint arXiv:2505.24357*.
- Dongjie Yang, Xiaodong Han, Yan Gao, Yao Hu, Shilin Zhang, and Hai Zhao. 2024a. PyramidInfer: Pyramid KV cache compression for high-throughput LLM inference. In *Findings of the Association for Computational Linguistics*, pages 3258–3270.

- June Yong Yang, Byeongwook Kim, Jeongin Bae, Beomseok Kwon, Gunho Park, Eunho Yang, Se Jung Kwon, and Dongsoo Lee. 2024b. No token left behind: Reliable KV cache compression via importance-aware mixed precision quantization. arXiv preprint arXiv:2402.18096.
- Peiyu Yang, Naveed Akhtar, Zeyi Wen, and Ajmal Mian. 2023a. Local path integration for attribution. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3173–3180.
- Peiyu Yang, Naveed Akhtar, Zeyi Wen, Mubarak Shah, and Ajmal Saeed Mian. 2023b. Re-calibrating feature attributions for model interpretation. In *International Conference on Learning Representations*.
- Yifei Yang, Zouying Cao, Qiguang Chen, Libo Qin, Dongjie Yang, Hai Zhao, and Zhi Chen. 2024c. KVSharer: Efficient inference via layerwise dissimilar KV cache sharing. arXiv preprint arXiv:2410.18517.
- Lu Ye, Ze Tao, Yong Huang, and Yang Li. 2024. ChunkAttention: Efficient self-attention with prefixaware KV cache and two-phase partition. In *Annual Meeting of the Association for Computational Linguistics*, pages 11608–11620.
- Zihao Ye, Lequn Chen, Ruihang Lai, Wuwei Lin, Yineng Zhang, Stephanie Wang, Tianqi Chen, Baris Kasikci, Vinod Grover, Arvind Krishnamurthy, and Luis Ceze. 2025. FlashInfer: Efficient and customizable attention engine for LLM inference serving. *Proceedings of Machine Learning and Systems*.
- Chengye Yu, Tianyu Wang, Zili Shao, Linjie Zhu, Xu Zhou, and Song Jiang. 2024a. TwinPilots: A new computing paradigm for GPU-CPU parallel LLM inference. In *ACM International Systems and Storage Conference*, pages 91–103.
- Hao Yu, Zelan Yang, Shen Li, Yong Li, and Jianxin Wu. 2024b. Effectively compress KV heads for LLM. *arXiv preprint arXiv:2406.07056*.
- Lingfan Yu, Jinkun Lin, and Jinyang Li. 2025. Stateful large language model serving with pensieve. In *Proceedings of the Twentieth European Conference on Computer Systems*, pages 144–158.
- Zhihang Yuan, Yuzhang Shang, Yang Zhou, Zhen Dong, Chenhao Xue, Bingzhe Wu, Zhikai Li, Qingyi Gu, Yong Jae Lee, Yan Yan, Beidi Chen, Guangyu Sun, and Kurt Keutzer. 2024. LLM inference unveiled: Survey and roofline model insights. *arXiv preprint arXiv:2402.16363*.
- Yuxuan Yue, Zhihang Yuan, Haojie Duanmu, Sifan Zhou, Jianlong Wu, and Liqiang Nie. 2024. WKVQuant: Quantizing weight and key/value cache for large language models gains more. *arXiv* preprint *arXiv*:2402.12065.

- Ahmet Caner Yüzügüler, Jiawei Zhuang, and Lukas Cavigelli. 2025. PRESERVE: Prefetching model weights and KV-cache in distributed LLM serving. *arXiv preprint arXiv:2501.08192*.
- Amir Zandieh, Majid Daliri, and Insu Han. 2025. QJL: 1-bit quantized JL transform for KV cache quantization with zero overhead. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 25805–25813.
- Amir Zandieh, Insu Han, Vahab Mirrokni, and Amin Karbasi. 2024. SubGen: Token generation in sublinear time and memory. *arXiv preprint arXiv:2402.06082*.
- Hailin Zhang, Xiaodong Ji, Yilin Chen, Fangcheng Fu,
 Xupeng Miao, Xiaonan Nie, Weipeng Chen, and Bin
 Cui. 2025a. PQCache: Product quantization-based
 KVCache for long context LLM inference. Proceedings of the ACM on Management of Data, 3(3):1–30.
- Minwei Zhang, Haifeng Sun, Jingyu Wang, Shaolong Li, Wanyi Ning, Qi Qi, Zirui Zhuang, and Jianxin Liao. 2025b. ClusterAttn: KV cache compression under intrinsic attention clustering. In Annual Meeting of the Association for Computational Linguistics, pages 14451–14473.
- Rongzhi Zhang, Kuang Wang, Liyuan Liu, Shuohang Wang, Hao Cheng, Chao Zhang, and Yelong Shen. 2024a. LoRC: Low-rank compression for LLMs KV cache with a progressive compression strategy. *arXiv* preprint arXiv:2410.03111.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. OPT: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068.
- Tianyi Zhang, Jonah Yi, Zhaozhuo Xu, and Anshumali Shrivastava. 2024b. KV cache is 1 bit per channel: Efficient large language model inference with coupled quantization. *Advances in Neural Information Processing Systems*, 37:3304–3331.
- Yanqi Zhang, Yuwei Hu, Runyuan Zhao, John Lui, and Haibo Chen. 2025c. Unifying KV cache compression for large language models with LeanKV. *arXiv* preprint arXiv:2412.03131.
- Yuxin Zhang, Yuxuan Du, Gen Luo, Yunshan Zhong, Zhenyu Zhang, Shiwei Liu, and Rongrong Ji. 2024c. CaM: Cache merging for memory-efficient LLMs inference. In *International Conference on Machine Learning*, pages 58840–58850. PMLR.
- Zeyu Zhang and Haiying Shen. 2025. FDC: Fast KV dimensionality compression for efficient LLM inference. *arXiv preprint arXiv:2408.04107*.

Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark W. Barrett, Zhangyang Wang, and Beidi Chen. 2023. H2O: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36:34661–34710.

Junqi Zhao, Zhijin Fang, Shu Li, Shaohui Yang, and Shichao He. 2024a. BUZZ: Beehive-structured sparse KV cache with segmented heavy hitters for efficient LLM inference. arXiv preprint arXiv:2410.23079.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Z. Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, and 3 others. 2023. A survey of large language models. *arXiv* preprint arXiv:2303.18223.

Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy, Tianqi Chen, and Baris Kasikci. 2024b. Atom: Lowbit quantization for efficient and accurate LLM serving. *Proceedings of Machine Learning and Systems*, 6:196–209.

Youpeng Zhao, Di Wu, and Jun Wang. 2024c. AL-ISA: Accelerating large language model inference via sparsity-aware KV caching. In *ACM/IEEE Annual International Symposium on Computer Architecture*, pages 1005–1017. IEEE.

Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph Gonzalez, Clark W. Barrett, and Ying Sheng. 2024. SGLang: Efficient execution of structured language model programs. *Advances in Neural Information Processing Systems*, 37:62557–62583.

Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu, Yibo Zhu, Xuanzhe Liu, Xin Jin, and Hao Zhang. 2024. DistServe: Disaggregating prefill and decoding for goodput-optimized large language model serving. In *USENIX Symposium on Operating Systems Design and Implementation*, pages 193–210.

Xiabin Zhou, Wenbin Wang, Minyan Zeng, Jiaxian Guo, Xuebo Liu, Li Shen, Min Zhang, and Liang Ding. 2024a. DynamicKV: Task-aware adaptive KV cache compression for long context LLMs. *arXiv preprint arXiv:2412.14838*.

Zixuan Zhou, Xuefei Ning, Ke Hong, Tianyu Fu, Jiaming Xu, Shiyao Li, Yuming Lou, Luning Wang, Zhihang Yuan, Xiuhong Li, Shengen Yan, Guohao Dai, Xiao-Ping Zhang, Yuhan Dong, and Yu Wang. 2024b. A survey on efficient inference for large language models. *arXiv preprint arXiv:2404.14294*.

A Design of Our Taxonomy

Our taxonomy follows a system behavior-oriented view of sKis introduced in § 2 and respects the domain boundary. Specifically, we classify techniques by their operational impact along three dimensions: temporal, spatial, and structural. This behavior-oriented perspective follows established practice in machine learning systems research (Xiao et al., 2018; Rajbhandari et al., 2020; Wang et al., 2021; Jiang et al., 2022; Qiu et al., 2024; Jiang et al., 2024) and aligns closely with how serving systems are actually built and optimized in practice, allowing diverse methods to be interpreted under a unified framework.

For example, many methods perform KV cache *selection* by identifying tokens (i.e., KV entries) that are more or less important for future computation. In our taxonomy, we do not treat selection itself as a category. In contrast, we classify methods based on the system action taken after selection:

- If unimportant KV entries are directly discarded to free memory, the method is categorized as KV cache eviction (cf. § 5.2.2) under KV representation and retention (structural dimension).
- If unimportant KV entries are offloaded to secondary storage for possible future retrieval and reload, the method falls under cross-device memory hierarchy (cf. § 4.1) in KV placement and migration (spatial dimension).
- If the tokens are retained in memory but excluded from computation, the method is considered token-level scheduling, which is categorized as KV-centric scheduling (cf. § 3.1) under KV execution and scheduling (temporal dimension).

In short, selection is treated as a preparatory step, not a classification criterion itself. This helps prevent ambiguity and ensures that each category in our taxonomy corresponds to a distinct system-level optimization behavior.

B Related Surveys

A rapidly growing body of surveys reviews efficient LLM inference and serving from algorithms to systems, e.g., speculative decoding (Xia et al., 2024), algorithms-to-systems serving (Miao et al., 2023), roofline-style analyses (Yuan et al., 2024), and broader resource-efficient or compression-oriented perspectives (Wang et al., 2024a; Bai et al., 2024; Li et al., 2024a; Xu et al., 2024b; Zhou et al., 2024b). In contrast, dedicated surveys that fo-

Table 4: Comparison with existing surveys related to efficient LLM inference or serving.

Survey	KV- centric	Serving only	No retrain	System metrics	Taxonomy principle
Xia et al. (2024)		1	1		Stages (draft, verify)
Miao et al. (2023)		✓	✓	✓	Layered pipeline (algorithm-, system-level)
Yuan et al. (2024)		✓	✓	✓	Layered pipeline (parameter-, algorithm-, system-, hardware-level)
Wang et al. (2024a)					Compression families (quantization, pruning, knowledge distillation, compact architecture design, dynamic networks)
Li et al. (2024a)		✓	✓	✓	System components (KV cache and memory, computation, cloud deployment, emerging research fields)
Bai et al. (2024)				✓	Lifecycle (architecture design, pre-training, fine-tuning, inference, system design)
Zhou et al. (2024b)				/	Layered pipeline (data-, model-, system-level)
Xu et al. (2024b)				✓	Layered pipeline (architecture, algorithm, systems)
Liu et al. (2025a)	1	✓	1		KV compression families (selective token, quantization, attention compression, hybrid)
Shi et al. (2024)	✓				Lifecycle (training, deploy, post-training)
Li et al. (2024b)	✓				Abstraction level (token-, model-, system-level)
This survey (sKis)	✓	✓	✓	✓	System behaviors (temporal, spatial, structural dimensions)

cus specifically on the KV cache remain rare. Shi et al. (2024) adopted a lifecycle-based taxonomy spanning training-stage, deploy-stage, and post-training optimizations. Li et al. (2024b) categorized KV cache management strategies into token-level, model-level, and system-level optimizations. Liu et al. (2025a) focused on compression strategies of the KV cache, such as selective token strategies, quantization, and attention compression.

Relation to existing surveys. While several recent surveys have covered related areas, this survey distinguishes itself by offering a distinct system behavior-oriented perspective. General LLM serving (inference) surveys typically treat KV cache optimization as a minor component within the broader pipelines. KV-specific surveys are closest to our topic. However, they mostly organize methods by lifecycle stages or by abstraction levels, leaving the serving-time system behavior of the KV cache largely unexamined. Instead, we concentrate exclusively on the sKis scope (i.e., serving-time, KV-centric, system metrics, no retraining or architecture change) and aim to provide a deeper understanding within this scope. By classifying methods according to their impact along temporal, spatial, and structural dimensions, our survey enables cross-behavior and behavior × objective analysis, which complements prior surveys and clarifies actionable research gaps for KV-centric serving. Table 4 shows a comparative summary.

C Supplementary Paper Categorization

Table 5 provides a supplementary mapping of all surveyed methods across the full taxonomy of

7 subcategories under 3 major optimization dimensions. The finer-grained categories in this table include (i) KV-centric scheduling (cf. § 3.1), (ii) pipelining and overlapping (cf. § 3.2), (iii) hardware-aware execution (cf. § 3.3), (iv) memory hierarchy KV orchestration (cf. § 4.1), (v) compute device KV orchestration (cf. § 4.2), (vi) KV cache compression (cf. § 5.1), and (vii) KV cache retention management (cf. § 5.2).

As discussed in § 2, to maintain structural clarity and prevent overly diffuse categorization, each method is primarily discussed under one or two key optimization categories that reflect its main contributions. These categories are denoted as primary category () in Tab. 5. However, some methods also touch upon additional optimization aspects that are not covered or elaborated in the main sections. For example, to support its "hardware-aware execution" design of decoupling prefill and decode phases across heterogeneous devices, Splitwise (Patel et al., 2024) incorporates a fine-grained layerwise transmission strategy that transmits the KV cache from the prefill node to the decode node and overlaps such KV cache transmission with the computation in the prefill phase. They serve as enabling mechanisms that make the decoupled strategy feasible and link Splitwise to the "device-level KV transfer" and "pipelining and overlapping" categories. We summarize these omitted associations in Tab. 5, denoted by \(\bar{\cupsto} \), to provide a more complete mapping for readers interested in cross-cutting techniques.

Venue Diversity. We can observe from the "Venue" column of Tab. 5 that the methods span a broad range of research communities. The pub-

Table 5: Full mapping of representative methods reviewed in this paper to their corresponding sKis categories. Methods are chronologically ordered with publication venues.

Methods	Venue				conomy of s	sKis		
$K_{V_{-\mathrm{ce}}}$	M Hai Pipelining and ntric schedulin	lemory l dware-a overlap g	Compunierarchy K aware execuping	te device K V orchestra tion	KV ca KV cac V orchestra ation	che retenti he compres tion	on manager Ssion	nent
FastServe (Wu et al., 2023)	arXiv		•		•	•		
SmoothQuant (Xiao et al., 2023)	ICML							
FlexGen (Sheng et al., 2023)	ICML							
CCompv1 (Park et al., 2023)	APSys							
vLLM (Kwon et al., 2023)	SOSP							
H ₂ O (Zhang et al., 2023)	NeurIPS							
Scissorhands (Liu et al., 2023)	NeurIPS							
TetriInfer (Hu et al., 2024b)	arXiv							
SubGen (Zandieh et al., 2024)	arXiv							
RoCo (Ren and Zhu, 2024)	arXiv							
WKVQuant (Yue et al., 2024)	arXiv							
MiKV (Yang et al., 2024b) FastDecode (He and Zhai, 2024)	arXiv arXiv							
	arXiv arXiv					•		
QAQ (Dong et al., 2024b) RAGCache (Jin et al., 2024)	arXiv arXiv		•					
AttAcc (Park et al., 2024)	ASPLOS							
FastGen (Ge et al., 2024)	ICLR							
StreamingLLM (Xiao et al., 2024b)	ICLR							
Preble (Srivatsa et al., 2024)	ICLR			•				
Keyformer (Adnan et al., 2024)	MLSys							
Atom (Zhao et al., 2024b)	MLSys							
PromptCache (Gim et al., 2024)	MLSys							
PyramidKV (Cai et al., 2024)	arXiv							
HeadKV (Yu et al., 2024b)	arXiv							
LESS (Dong et al., 2024a)	arXiv							
D ₂ O (Wan et al., 2024)	arXiv							
Splitwise (Patel et al., 2024)	ISCA		•					
ALISA (Zhao et al., 2024c)	ISCA						•	•
Infinite-LLM (Lin et al., 2024)	arXiv		•					
Ada-KV (Feng et al., 2024)	arXiv							
Mooncake (Qin et al., 2024)	arXiv							
DistServe (Zhong et al., 2024)	OSDI							
InfiniGen (Lee et al., 2024)	OSDI							
CachedAttention (Gao et al., 2024)	ATC							•
LazyLLM (Fu et al., 2024)	arXiv							
KVMerger (Wang et al., 2024b)	arXiv							
vTensor (Xu et al., 2024a)	arXiv							
KIVI (Liu et al., 2024f)	ICML							
CHAI (Agarwal et al., 2024)	ICML							
CaM (Zhang et al., 2024c)	ICML			_				
MuxServe (Duan et al., 2024)	ICML					•		•
Quest (Tang et al., 2024)	ICML							
SparQAttention (Ribar et al., 2024)	ICML					_		
DéjàVu (Strati et al., 2024)	ICML					0		•
CacheGen (Liu et al., 2024e)	SIGCOMM							
DecoQuant (Liu et al., 2024d)	ACL							
NACL (Chen et al., 2024)	ACL							
PyramidInfer (Yang et al., 2024a)	ACL							
ChunkAttention (Ye et al., 2024) InstInfer (Pan et al., 2024)	ACL orViv							
, ,	arXiv SYSTOR		0					
TwinPilots (Yu et al., 2024a)	arXiv		•					
GEAR (Kang et al., 2024) LoRC (Zhang et al., 2024a)	arXiv arXiv							
Lorce (Zhang et al., 2024a)	ai /XIV							

Continued on next page

Methods Venue Taxonomy of sKis Memory hierarchy KV orchestration management Memory hierarchy hierarchy kV orchestration management Memory hierarchy Pipelining and overlapping KV-centric scheduling LayerKV (Xiong et al., 2024) arXiv 0 0 CCompv2 (Park and Egger, 2024) **PACT** C KVSharer (Yang et al., 2024c) arXiv LAMPS (Shahout et al., 2024) arXiv • 0 BUZZ (Zhao et al., 2024a) arXiv Palu (Chang et al., 2025b) arXiv TokenSelect (Wu et al., 2024b) arXiv 0 LoongServe (Wu et al., 2024a) SOSP • € SKVQ (Duanmu et al., 2024) arXiv TOVA (Oren et al., 2024) **EMNLP** VATP (Guo et al., 2024) **EMNLP** L2KV (Devoto et al., 2024) **EMNLP** FastSwitch (Shen et al., 2024) arXiv 0 KVQuant (Hooper et al., 2024) NeurIPS CQ (Zhang et al., 2024b) NeurIPS ZipCache (He et al., 2024) NeurIPS SnapKV (Li et al., 2024c) NeurIPS MiniCache (Liu et al., 2024a) NeurIPS InfLLM (Xiao et al., 2024a) NeurIPS Radix Attention (Zheng et al., 2024) NeurIPS Loki (Singhania et al., 2024) NeurIPS DynamicKV (Zhou et al., 2024a) arXiv MemServe (Hu et al., 2024a) arXiv RetrievalAttention (Liu et al., 2024b) arXiv QJL (Zandieh et al., 2025) AAAI **HPCA** VQ-LLM (Liu et al., 2025b) xKV (Chang et al., 2025a) arXiv SQuat (Wang et al., 2025) arXiv vAttention (Prabhu et al., 2025) **ASPLOS** C PAPI (He et al., 2025) **ASPLOS** Pensieve (Yu et al., 2025) EuroSys • C LeanKV (Zhang et al., 2025c) arXiv AsyncKV (Dong et al., 2025) arXiv ShadowKV (Sun et al., 2024) arXiv C CAKE (Qin et al., 2025) **ICLR** ThinK (Xu et al., 2025b) **ICLR** MatryoshkaKV (Lin et al., 2025a) **ICLR** MagicPIG (Chen et al., 2025b) **ICLR** MLSys QoQ (Lin et al., 2025b) FlashInfer (Ye et al., 2025) MLSys • MLSys Neo (Jiang et al., 2025b) • NSNQuant (Son et al., 2025) arXiv PRESERVE (Yüzügüler et al., 2025) arXiv C ReCalKV (Yan et al., 2025) arXiv FDC (Zhang and Shen, 2025) arXiv ClusterKV (Liu et al., 2024c) DAC PQCache (Zhang et al., 2025a) **SIGMOD** C SepLLM (Chen et al., 2025a) **ICML** CommVQ (Li et al., 2025) **ICML** LaCache (Shi et al., 2025) **ICML** ClusterAttn (Zhang et al., 2025b) ACL RefreshKV (Xu et al., 2025a) ACL OTT (Su et al., 2025) ACL KVPR (Jiang et al., 2025a) ACL

^{• =} primary category with main analysis; • = secondary category omitted or only briefly mentioned in our paper to maintain focused classification.

lication venues include top-tier machine learning and artificial intelligence conferences (e.g., ICLR, ICML, NeurIPS, AAAI), natural language processing venues (e.g., ACL, EMNLP), systems and architecture conferences (e.g., ASPLOS, ISCA, HPCA, ATC, EuroSys, OSDI, SOSP, SC, SIGCOMM, DAC), and interdisciplinary forums such as ML-Sys and SIGMOD. We also include some impactful arXiv submissions. This diversity underscores the inherently cross-cutting nature of KV cache optimization, which lies at the intersection of model serving and system efficiency. It also highlights the growing recognition of this topic across various research communities.

D Key Takeaways

Through a comprehensive literature review of sKis, we have discovered takeaways across several domains. These include scheduling & overlapping, hardware-aware execution, placement & migration, compression, and eviction.

D.1 Scheduling and Overlapping

KVS and OVLP directly target runtime stalls. KVS prioritizes limited resources for the most reusable and latency-sensitive work; OVLP is also a type of scheduling, aligning compute with data transfer to fill pipeline bubbles.

ATakeaway:

- ✓ KVS is a multi-objective optimization problem. Modern schedulers often prioritize KV usage over time rather than FLOPS, and KV reusedriven scheduling is the default paradigm.
- ✓ KVS is enhanced by prediction. Lightweight predictors plus a robust policy outperform traditional FCFS or SJF schemes (Hu et al., 2024b; Qin et al., 2024; Shahout et al., 2024).
- ✓ The key to OVLP is to perform at the true bottleneck with asymmetric pipelines. For example, keep compute-bound prefill on GPU, and overlap memory-bound decode attention and KV with I/O or collective communication.
- ✓ Preferring recompute to transfer, e.g., partially recomputing KV while streaming the rest (Jiang et al., 2025a), or prefetching KV caches into L2 during collectives (Dong et al., 2025; Yüzügüler et al., 2025), can substantially reduce pipeline bubbles, especially when bandwidth is the bottleneck.

D.2 Hardware-aware Execution

HAE improves throughput, reduces mean/tail latency, and extends servable context without retraining, by decoupling phases and mapping execution to hardware capabilities.

QTakeaway:

- ✓ Compute should follow hardware capabilities. When executing on a given device, it is critical to specialize kernels, tiling, and memory layouts to that device.
- ✓ Create KV locality within the device rather than moving KV across devices. It is effective to keep hot KV caches close to the compute.
- ✓ Compute-intensive prefill and memory-bound decode benefit from phase-specific execution mappings (cf. § 3.3.2).
- ✓ HAE should adapt to the access granularity and parallelism of the target device.

D.3 Placement and Migration

MHO and CDO govern where KV caches reside across the memory hierarchy and how they transfer during serving. They act directly on interconnect bandwidth bottlenecks, with GPU memory relief emerging as a by-product of tiering and offloading.

QTakeaway:

- ✓ It is a common MHO pattern to keep only futureuseful KV caches on the GPU, demote the rest to CPU or SSD, and reload guided by attention cues. Cost models can be effectively used to choose CPU, GPU, SSD paths (Sheng et al., 2023; Jin et al., 2024).
- ✓ Most MHO and CDO solutions overlap I/O transfers with compute or collectives to hide latency, although they often serve OVLP as a secondary category.
- ✓ Under interconnect bottlenecks, co-adaptation of transfer paths, precisions, or decoding strategies can reduce TTFT and SLO violations compared with static schemes (Zhong et al., 2024; Liu et al., 2024e; Shen et al., 2024).
- ✓ Migration granularity and path should align with attention access patterns and device access units.
- ✓ Prefetch-evict co-optimization remains rare. The field would benefit from a unified objective that jointly accounts for prefetch deadlines and eviction risk.

D.4 KV Cache Compression

KV caches can quickly overwhelm the memory capacity of GPUs and pose bandwidth pressure as context length or batch size increases, since the size of the KV cache scales linearly with these two factors. Consequently, prior works have proposed various approaches to directly compress the KV cache, such as quantization, low-rank approximation, and structural compression.

QTakeaway:

- ✓ Outlier handling dominates performance at low bitwidths or ranks (Su et al., 2025). Isolating outliers (e.g., higher bitwidths) for value-level compression methods prevents worst-case error explosions.
- ✓ Recent advances trend toward applying vector quantization (VQ) for KV cache quantization, and they often reach very low-bit (i.e., 1-2 bits) quantization with modest quality loss (Zhang et al., 2024b; Liu et al., 2025b; Son et al., 2025; Li et al., 2025). VQ (Gray, 1984) is a popular technique to represent high-dimensional data using a smaller set of representative vectors, known as codebooks. Variants of VQ like product quantization (Jegou et al., 2010) and additive quantization (Babenko and Lempitsky, 2014) have been proposed to applied to KVCC.
- ✓ KVCC has been developed mostly at the algorithm level, while system-level integration is thin (cf. Fig. 6). Thus, memory reductions often fail to translate into lower mean/tail latency or higher throughput unless KVCC is co-designed with execution, migration, and runtime control.

We further discuss the co-design of KVCC with execution, migration, and runtime control (the last akeaway) as follows: (i) Co-design with execution: quantization/de-quantization and low-rank updates can be fused into attention kernels or overlapped with compute, so compression overhead does not re-introduce stalls in the decode pipeline; (ii) Co-design with migration: aligning compressed packing units with device access units ensures that memory footprint reductions translate into fewer, fully utilized transfer chunks that fit overlap windows. (iii) Co-design with runtime control: exposing tunable parameters (e.g., bitwidth, rank, sparsity) to the runtime and adjusting them under SLOs remains an opportunity beyond static configurations.

D.5 KV Cache Eviction

KV cache eviction decides which past tokens remain resident under tight memory and bandwidth budgets, so that long contexts can be served. It operates in both phases and trades memory and transfer cost against utility to the attention compute.

ATakeaway:

- ✓ KV cache eviction is important in both prefill and decode. The former focuses on the KV cache to be computed, while the latter focuses on the KV cache that has been computed.
- ✓ Most systems retain a small recent window, a tiny set of "attention sink" anchor tokens (Xiao et al., 2024b; Gu et al., 2025), and a few "heavy hitters" (Zhang et al., 2023) identified by cumulative attention.
- ✓ Token importance should not be judged by attention scores alone. KV norms provide strong and low-overhead signals (Guo et al., 2024; Devoto et al., 2024). We also recommend calibrating token importance scores before using them for eviction (Sundararajan et al., 2017; Smilkov et al., 2017; Yang et al., 2023a,b).
- ✓ It is effective to use heterogeneous budgets across layers or heads, rather than a uniform upper bound (Cai et al., 2024; Yang et al., 2024a; Wan et al., 2024; Qin et al., 2025; Shi et al., 2025; Feng et al., 2024). For example, shallow layers often deserve larger retention, while deeper layers emphasize global semantics and tolerate more sparsity.
- ✓ Pairing KV cache eviction with similarity-based recall or merge is stronger than hard deletion, preserving salient context under tight budgets and improving long-context consistency.

E Behavior-behavior Synergy Compute

As discussed in § 6, Fig. 6 presents a behavior-behavior synergy network that summarizes how often behaviors co-occur within the same paper across seven behaviors, including KVS, OVLP, HAE, MHO, CDO, KVCC, and KVRM. Below we detail the procedure for calculating the normalized co-occurrence strengths for behavior pairs, which are reflected by the edge thicknesses in Fig. 6.

Let $\mathcal{B} = \{\text{KVS, OVLP, HAE, MHO, CDO, KVCC, KVRM}\}$ denote the set of system behaviors and \mathcal{P} the set of papers. For paper $p \in \mathcal{P}$ and behavior $i \in \mathcal{B}$, let the categorical label be $\ell_{p,i} \in \{\text{P,S,NA}\}$, which means primary category

	KVS	OVLP	HAE	МНО	CDO	KVCC	KVRM
KVS	_	2.25	4.75	3.5	1.75	0	6
OVLP	2.25	_	9	6.5	6	2.75	2
HAE	4.75	9	_	3.25	9.5	1.25	3
MHO	3.5	6.5	3.25	_	2	3	5.25
CDO	1.75	6	9.5	2	_	1.25	2.25
KVCC	0	2.75	1.25	3	1.25	_	1.75
KVRM	6	2	3	5.25	2.25	1.75	_

Table 7: Raw (pre-normalization) co-occurrence matrix that encodes the weighted co-occurrence strength between system behaviors across papers.

(•), secondary category (•), or no category. Each $\ell_{p,i}$ can be observed from Tab. 5. We map labels to numeric weights by $\omega(\ell_{p,i}) = \mathbb{1}_{[\ell_{p,i}=\mathsf{P}]} + \alpha \mathbb{1}_{[\ell_{p,i}=\mathsf{S}]}$ with $\alpha=0.5$, where $\mathbb{1}_{[\cdot]}$ is the indicator function that equals 1 when the stated condition holds and 0 otherwise.

Constructing raw co-occurrence. The raw co-occurrence matrix $C \in \mathbb{R}^{|\mathcal{B}| \times |\mathcal{B}|}$ aggregates pairwise co-appearance strength, as shown in Tab. 7. Each cell C_{ij} of the behavior pair i,j is defined by summing the per-paper products of their weights:

$$C_{ij} = \sum_{p \in \mathcal{P}} \omega(\ell_{p,i})\omega(\ell_{p,j}).$$

Equivalently, each paper p contributes 1 if both behaviors are primary, α if one is primary and the other secondary, α^2 if both are secondary, and 0 otherwise. The matrix C is symmetric.

Constructing normalized synergy. While the raw co-occurrence matrix C captures absolute overlap, it is biased by marginal popularity, because the behaviors with larger research density tend to have larger C_{ij} even without specific affinity. We therefore normalize C using the Tanimoto coefficient. We define the per-behavior squared weight Q_i :

$$Q_i = \sum_{p \in \mathcal{P}} w_{p,i}^2.$$

Then the Tanimoto-normalized synergy matrix $S \in \mathbb{R}^{|\mathcal{B}| \times |\mathcal{B}|}$ reflects relative co-occurrence strength on a [0,1] scale, as shown in Tab. 8. Each cell in S_{ij} of the behavior pair i,j is defined as the ratio of their shared weighted presence to their squared union:

$$S_{ij} = \frac{C_{ij}}{Q_i + Q_j - C_{ij}}.$$

Compared to C_{ij} , this score controls marginal sizes and is visualized in Fig. 6. We draw an undirected edge between behaviors i and j iff $S_{ij} > \theta$, where

	KVS	OVLP	HAE	MHO	CDO	KVCC	KVRM
KVS	_	0.09	0.17	0.12	0.08	0	0.14
OVLP	0.09	_	0.43	0.28	0.39	0.06	0.05
HAE	0.17	0.43	_	0.10	0.53	0.02	0.06
MHO	0.12	0.28	0.10	_	0.08	0.06	0.11
CDO	0.08	0.39	0.53	0.08	_	0.03	0.05
KVCC	0	0.06	0.02	0.06	0.03	_	0.02
KVRM	0.14	0.05	0.06	0.11	0.05	0.02	

Table 8: Normalized synergy matrix that encodes relative co-occurrence strength between behaviors across papers. Scores greater than the threshold $\theta=0.14$ are highlighted and visualize in Fig. 6 accordingly.

we set the threshold $\theta = 0.14$; edges below the threshold are omitted to reduce clutter. Edge thickness is proportional to S_{ij} .

F Benchmarking for sKis

In this section, we focus on system-performance benchmarking during serving, which measures actual performance metrics like latency, throughput, service-level objectives (SLOs), KV cache memory and bandwidth, and energy. In contrast, task or quality benchmarks focus on datasets and accuracy metrics. In our survey they serve only as quality gates and are not primary evaluation objectives. We refer interested readers to another survey (Li et al., 2024b) for details.

Many popular inference frameworks or systems, such as vLLM (Kwon et al., 2023), TensorRT-LLM (NVIDIA, 2023), and DeepSpeed-Inference (Aminabadi et al., 2022), provide benchmark scripts that measure system metrics for their local checks but remain framework-specific. We therefore view them as systems under test rather than the benchmark itself. In this section, we survey benchmarking efforts that provide a platform-and framework-agnostic way to obtain system measurements. They primarily fall into two categories: client-side tools and benchmark suites.

Client-side tools define and enforce metric semantics. Using one tool across systems yields directly comparable numbers. LLMPerf (Ray, 2024) targets API benchmarking and provides system metric measurement on service endpoints. NVIDIA NIM benchmarking guide (NVIDIA, 2025b) defines the common metrics of time to first token (TTFT), end-to-end request latency, intertoken latency (ITL), tokens per second (TPS), and requests per second (RPS). The companion tool GenAI-Perf (NVIDIA, 2025a) emits the defined metrics and implements the stable-window analysis

across OpenAI-API-compatible backends. However, although client-side tools offer specific metrics for LLM-based applications, we find inconsistent metric definitions and measurements across different tools.

Benchmark suites mean standardized packages of workloads, procedures, and reporting rules that specify what to run, how to run it, and what to report. Such suites typically cover multiple systems and hardware and enable reproducible comparisons. MLPerf Inference (Reddi et al., 2020) emphasizes inference system comparison, and its v5.0 includes LLM scenarios with accuracy validation. LLM-Inference-Bench (Chitty-Venkata et al., 2024) evaluates the inference performance of the LLaMA model family across a variety of hardware platforms. BALI (Jurkschat et al., 2025) measures LLM inference across six frameworks or acceleration approaches. It divides inference into three measured stages: setup, tokenize, and generate, and supports two settings: a technical setting with a fixed number of tokens, and a prompt-to-answer setting that includes tokenization.

G The Use of AI assistants

We used ChatGPT minimally for wording and grammar suggestions. No technical claims, taxonomy decisions, or analyses were produced by the assistant. All content was authored, verified, and edited by the authors.