

# SPECMD: A COMPREHENSIVE STUDY ON SPECULATIVE EXPERT PREFETCHING

Duc Hoang\*    Ajay Jaiswal    Mohammad Samragh    Minsik Cho  
Apple

## ABSTRACT

Mixture-of-Experts (MoE) models enable sparse expert activation, meaning that only a subset of the model’s parameters is used during each inference. However, to translate this sparsity into practical performance, an expert caching mechanism is required. Previous works have proposed hardware-centric caching policies, but how these various caching policies interact with each other and different hardware specification remains poorly understood. To address this gap, we develop **SpecMD**, a standardized framework for benchmarking ad-hoc cache policies on various hardware configurations. Using SpecMD, we perform an exhaustive benchmarking of several MoE caching strategies, reproducing and extending prior approaches in controlled settings with realistic constraints. Our experiments reveal that MoE expert access is not consistent with temporal locality assumptions (e.g LRU, LFU). Motivated by this observation, we propose **Least-Stale**, a novel eviction policy that exploits MoE’s predictable expert access patterns to reduce collision misses by up to  $85\times$  over LRU. With such gains, we achieve over 88% hit rates with up to 34.7% Time-to-first-token (TTFT) reduction on OLMoE at only 5% or 0.6GB of VRAM cache capacity.

## 1 INTRODUCTION

Mixture-of-Experts (MoE) architectures scale language models to unprecedented sizes through sparse computation. By activating only  $k$  of  $N$  experts per token, MoE models like Mixtral-8x7B Jiang et al. (2024), DeepSeek-V2 DeepSeek-AI et al. (2024), and OLMoE Muennighoff et al. (2025) achieve sub-linear computational cost while maintaining large model capacity. However, this efficiency comes at a cost: a naive deployment of MoE models requires all  $N$  experts to remain memory-resident, making deployment prohibitively expensive. For instance, assuming 16 bits per weight, OLMoE-1B-7B activates 2.39GB yet needs 12.0GB ready in memory. Mixtral-8x7B requires 88GB storage yet only uses 24GB. This linear memory scaling limits the MoE deployment on memory-constrained devices.

Instead of keeping the full MoE model in active memory, an alternative is to store only the subset of patronized experts in a fast but capacity-limited cache, making memory management central to performance which can be influenced by the four aspects: **•routing**: how to select experts to be cached, **•prefetching** how to predict and cache the experts needed in future, **•eviction** how to free space in the cache for demanded experts, and **•miss handling** how to respond when required experts are absent from the cache. The importance of the above four aspects varies depending on hardware specifications. For instance, routing is capacity-dependent, prefetching is constrained by bandwidth, eviction matters in low-capacity settings, and miss handling is critical in balancing user-experience and task-performance.

Existing approaches tackle the cache management challenges through routing modification Skliar et al. (2025), LRU caching Eliseev & Mazur (2023), and dynamic precision Tang et al. (2024); Zhu et al. (2025), but their policies are heavily hardware-centric. Hence, they overfit to very specific set of circumstances, leaving general insights into policy interactions underexplored. This raises a fundamental research question: *when, where, and which combinations of these policies would offer the most benefits across different deployment scenarios?* It is for this end that we develop **SpecMD** (Speculative MoE Decoding), a standardized benchmarking framework designed

\*Corresponding author: dhoang2@apple.com

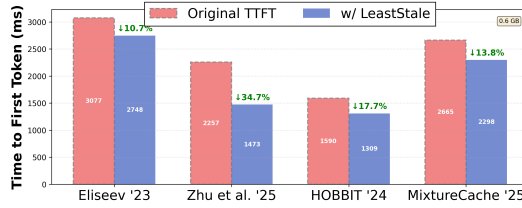


Figure 1: Least-Stale eviction as drop-in improvement across baseline approaches. Our novel eviction policy achieves 10.7–34.7% TTFT reduction on OLMoE-1B-7B at 5% cache capacity (0.6 GB).

to enable reproducible comparison of MoE caching strategies. Through this framework, we gain three capabilities: (1) controlled comparison of prior approaches under identical conditions, (2) exploration of policy interactions and compositions, and (3) characterization of performance across the full hardware constraint space.

Using SpecMD, we conducted an exhaustive benchmarking study of MoE caching policies. We tested LRU caching Eliseev & Mazur (2023), fetch cascade Tang et al. (2024), and expert substitution Zhu et al. (2025) under controlled and comparable conditions. We evaluate four MoE models (Mixtral-8x7B, OLMoE-1B-7B, Phi-3.5-MoE, Qwen1.5-MoE) and three different hardware configurations. From our studies, we made the following observations:

- **MoE’s expert access can cause wrongful eviction.** In MoE, experts are accessed following a deterministic sequential layer pattern, not recency-based reuse, rendering LRU and LFU policies perform quite poorly.
- **Dynamic prefetching is better for cache.** We observe dynamic, score-based prefetching outperforms top-k approaches; despite lower overall prediction accuracy, score-based prefetching yields higher hit-rates and makes use of available bandwidth more effectively.
- **Policy rankings shift across hardware regimes.** As expected, the best-performing strategy at 1% capacity differs from that at 5% or 25%; bandwidth-limited and capacity-limited scenarios favor different approaches. No single set of policy dominates, but there are trends (see Figure 9 in Appendix).

From these findings, we propose **Least-Stale**, an eviction policy that balances temporal and spatial awareness to minimize cache collision misses. Our contributions are:

- **A Framework for Policy Research.** We introduce an open-source policy-first framework for MoE Cache to accelerate policy prototyping and interaction.
- **Policies studies and Benchmarking.** Using SpecMD, we study policies interaction on different models and constraints and reveal key-insights for policy-makers.
- **Least-Stale Eviction Policy.** Building on our findings, we propose Least-Stale, a flexible eviction policy that greatly reduces collision-misses by up to  $85\times$  versus LRU at 1% capacity, achieving 88-92% hit rates. We show in Figure 1 Least-Stale provides consistent gains across diverse baseline strategies.

## 2 RELATED WORK

MoE memory constraints are addressed through several complementary approaches: routing modification, expert caching with prefetching and quantization. We focus on expert caching approaches, as they directly target cache management, i.e., how to decide which experts to load or evict and how to handle cache misses. Existing work evaluates 1-2 policy combinations in isolation, leaving policy interactions and performance across hardware constraints largely uncharacterized.

**Prefetching.** Eliseev & Mazur (2023) introduce speculative prefetching for memory-constrained MoE inference. Their approach predicts next-layer expert requirements using current-layer hidden states, enabling single-layer-ahead speculation. HOBBIT Tang et al. (2024) extends this by leveraging gating input similarity across adjacent layers to prefetch multiple layers ahead. Both approaches predict based on architectural proximity but do not account for expert importance when bandwidth is limited. Zhu et al. (2025) addresses this through score-based prefetching, which filters experts by

gate score percentile rather than fixed top- $k$  counts. This prioritizes high-importance experts when bandwidth cannot accommodate all predictions.

**Eviction.** Eviction policies determine which cached experts to remove when capacity is exceeded. Eliseev & Mazur (2023) employ LRU caching, keeping a fixed number of recently-used experts across all layers. Tang et al. (2024) proposes LHU (Least High-precision Used), a variant of LFU that prioritizes retaining high-precision experts. Zhu et al. (2025) use score-based eviction, removing experts with the lowest accumulated activation scores.

**Miss handling.** For cache misses, Eliseev & Mazur (2023) trigger blocking loads from CPU memory. HOBBIT introduces precision cascading: dynamically loading lower-precision versions (int4 for float16, int2 for int8) based on gating score thresholds Tang et al. (2024). Zhu et al. (2025) employ expert substitution, replacing cache-miss experts with functionally similar cached experts whose scores fall within tolerance thresholds.

**Routing.** MixtureCache Skliar et al. (2025) takes a fundamentally different approach by modifying routing rather than optimizing cache policies. They bias expert selection toward already-cached experts through learnable parameters applied to router logits. Their system employs LRU eviction but no prefetching mechanism, instead treating cache misses as inevitable and modifying routing to reduce misses.

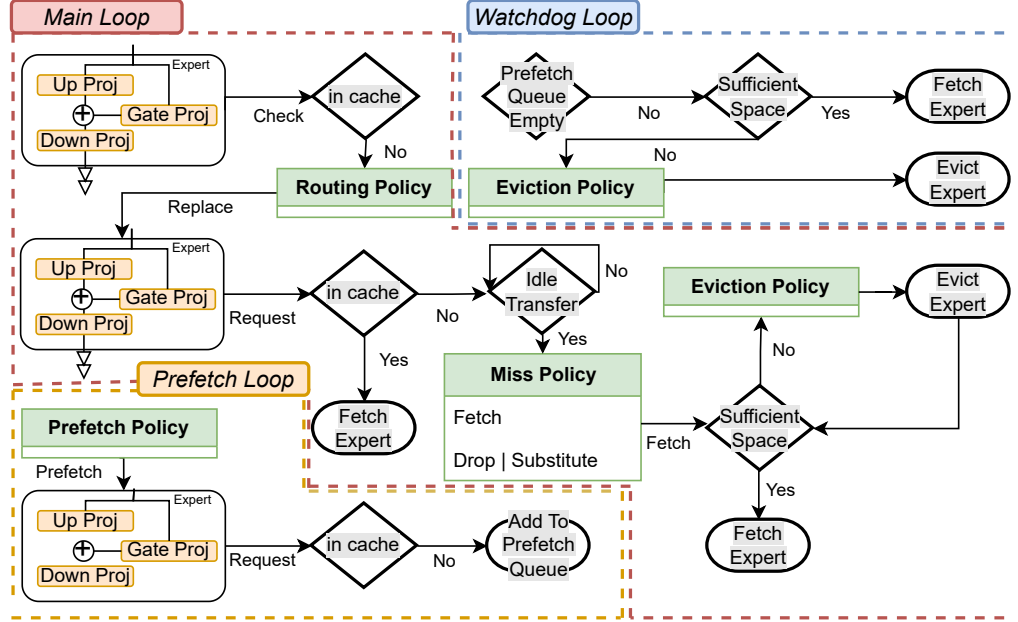


Figure 2: Main(red) and Prefetch (yellow) loops operate synchronously and are directly invoked by the model’s after gate. Watchdog (blue) operates asynchronously, constantly monitoring the prefetch queue.

### 3 SPECMD: FRAMEWORK AND POLICY SPACE

We present a framework that prioritizes ease of use (straightforward integration with HuggingFace models), accessibility (no specialized hardware required), and policy modularity. We aim to enable straightforward policy comparison across models and hardware on consumer-grade GPUs\*.

#### 3.1 POLICY DESIGN SPACE

Cache policies interact across four dimensions, shown in Figure 2, each affecting *quality*, *speed*, or both. Understanding these interactions is the central challenge. We systematically explore this

\*open-sourcing SpecMD is underway.

design space to characterize policy behaviors and identify optimal configurations under different quality-memory constraints.

### 3.1.1 ROUTING POLICIES

When an expert is requested by the gate module in an MoE block (top left of Figure 2), it may not exist in the cache. In this case, the routing policy is the first module that we may visit. Routing policies (shown as a green policy box in the Main Loop of Figure 2) determine whether to influence expert selection based on cache state, impacting both *performance* and *speed*. This represents a trade-off between cache efficiency and routing integrity.

**Standard routing** preserves original model behavior using unmodified router logits, maintaining 100% routing fidelity. All other policy dimensions operate around the model’s original expert selections.

**Cache-aware routing** biases selection toward cached experts by boosting GPU-resident expert logits Skliar et al. (2025). Let  $z = G(x)$  represent original router logits,  $C$  denote cached expert indices, and  $\mathbb{1}_C \in \{0, 1\}^{|E|}$  indicates cached experts. Modified logits are

$$z' = z + \lambda \cdot \Delta_{avg} \cdot \mathbb{1}_C, \quad (1)$$

where  $\lambda$  is hyperparameter used for tuning and  $\Delta_{avg}$  is the historical average logit value. Modified logits  $z'$  pass through softmax and top- $k$  selection as usual.

### 3.1.2 CACHE MISS POLICIES

After the routing policy is applied, the expert may still not be available in cache. Cache miss policies, as the name suggested, handle missing experts from cache (this is showed in the Main Loop of Figure 2). When cache misses occur, these policies can force immediate eviction to make space and may stall the model until the needed experts are fetched, substituted, or dropped. Unlike eviction/prefetching (primarily affecting speed), miss handling impacts both **quality** and **speed**. We implement three cache miss policies:

**Fetch policies** always fetch missing experts from CPU memory, preserving routing integrity at the cost of synchronous transfer latency. Basic Fetch guarantees 100% fidelity but degrades speed when misses are frequent. We introduce two precision-aware variants: *Fetch-lowest-bit* (e.g., int8:4) fetches the lowest available precision to minimize latency, while *Fetch-priority* (e.g., int8:4:2) always fetches highest precision with fallback levels (int8→int4→int2) for graceful degradation under capacity constraints.

**Drop policies** selectively skip cache-miss experts. Skliar et al. (2025) show MoE models are robust to expert variations: dropping top-ranked experts compromises performance, but models are resilient beyond rank 2, especially granular MoEs with many active experts. They demonstrate minimal perplexity degradation when swapping rank 3+ experts, which we leverage similarly by dropping only low-rank experts and always fetching high-rank ones.

**Substitution policies** replace cache-miss experts with functionally similar cached experts based on score proximity Zhu et al. (2025), but quality degrades significantly in practice (Section 5.3).

### 3.1.3 EVICTION POLICIES

After the cache-miss policy is applied, an expert may need to be loaded into the cache. Eviction policies determine which experts to remove when cache capacity is exceeded, and play a fundamental role in both the Main and Watchdog loop as seen Figure 2. In isolation, eviction policies influence **speed** aspect of a caching system. We support a couple broad categories of eviction policies (full details is in Appendix B):

- **Temporal policies.** LRU (Least Frequently Used) and LFU (Least Frequently Used) are classic examples that track and select the victim by access-time or access-frequency.
- **Spatial policies.** SB (Score-Based) and FLD (Furthest Layer Distance) select the victim based on its physical characteristic such as gate’s logits or distance from current layer.

**Algorithm 1** Gate Hook Processing

---

```

1: Input: Gate outputs, layer  $l$ 
2:  $(\text{experts}, \text{weights}, \text{logits}, \text{router\_weights}) \leftarrow$  Gate outputs
3: // Phase 1: Routing modification
4:  $\text{logits}' \leftarrow \text{RoutingPolicy.apply}(l, \text{logits})$ 
5:  $\text{experts}' \leftarrow \text{top\_k}(\text{softmax}(\text{logits}'))$ 
6:  $\text{weights}' \leftarrow \text{router\_weights}[\text{experts}']$ 
7: // Phase 2: Cache management (during expert execution)
8: for each expert  $e$  in  $\text{experts}'$  do
9:   if  $e \notin \text{GPU\_Cache}$  then
10:      $\text{MissPolicy.handle}(e)$ 
11:   end if
12: end for
13: // Phase 3: Prefetch submission
14: if  $l$  has prefetch targets  $\mathcal{T}$  then
15:   for each target layer  $t \in \mathcal{T}$  do
16:      $\text{predicted} \leftarrow \text{PrefetchPolicy.predict}(t)$ 
17:      $\text{PrefetchQueue.submit}(\text{predicted}, t)$ 
18:   end for
19: end if
20: Return:  $(\text{experts}', \text{weights}', \text{logits}', \text{router\_weights})$ 

```

---

## 3.1.4 PREFETCH STRATEGIES

Prefetch strategies (illustrated in the Prefetch Loop of Figure 2) load experts before they are requested, primarily influencing **speed**. The idea is to utilize the gate from an upcoming layer together with the activations of the current layer to predict future required experts. In our experiments we resort to prefetching only the very next layer experts.

Prefetch operates asynchronously via the Watchdog Loop (shown in blue in Figure 2) that monitors capacity and enqueues prefetch requests to the GPU when space is available.

We support two configuration dimensions in SpecMD:

- **Top-k prefetch with overfetch.** The standard approach prefetches the top-k experts per layer based on router weight magnitude, where  $k$  is model-specific ( $k = 8$  for OLMoE,  $k = 2$  for Mixtral). We introduce an overfetch factor that multiplies the prefetch count to increase bandwidth utilization: with overfetch factor 1.5, OLMoE prefetches 12 experts instead of 8. This trades prefetch precision for cached experts recall at the expense of bandwidth.
- **Score-based approach.** Following Zhu et al. (2025), this strategy filters experts by softmax score rather than fixed top-k count and prefetch all experts above a configurable percentile threshold (the default is 80th percentile). Unlike top-k approaches, this method adapts the number of prefetched experts based on score distribution.

## 3.2 IMPLEMENTATION OVERVIEW

To avoid modifying a model’s internal implementation, we utilized Pytorch forward hooks to intercepts and modify MoE gate outputs. We also standardize the gate-output protocol interface across architectures (Mixtral, OLMoE, Qwen) for more effective code-reuse. In Algorithm 1 we show the core hook processing across three phases: routing modification, cache management, and prefetch submission. Complete implementation details appear in Appendix A.

## 4 LEAST-STALE EVICTION POLICY

In this section, we present our novel policy contribution Least-Stale, which addresses some of the gaps presented in traditional eviction policies for performance improvement.

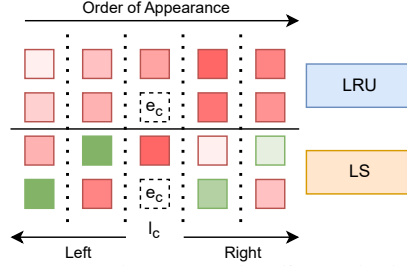


Figure 3: **Top:** LRU eviction as a purely temporal policy. Blocks with higher intensity indicate higher eviction priority; consequently, experts from the immediately following layer are most likely to be evicted. **Bottom:** Illustration of the Least-Stale (LS) policy, which combines temporal and positional signals. Cached experts are marked as current (green, accessed in the current forward pass) or stale (red, accessed in a previous pass), while layer index encodes position. Higher color intensity corresponds to higher eviction probability. Stale experts have a higher eviction priority than current ones.

#### 4.1 MOTIVATION

MoE expert access follows a deterministic front-to-back layer sequence (layer  $0 \rightarrow 1 \rightarrow 2 \rightarrow \dots \rightarrow 15$ ) within each forward pass. Once we move on from a layer, that layer’s experts will not be needed until the next forward pass. Temporal-only policies such as LRU and LFU utilize access-time and access-frequency as indicators for expert importance, and thus treat a recently accessed expert as important even though it is now irrelevant to the current forward pass. On the other hand, spatial-only policies like FLD and Score-Based use layer distance or router scores to indicate importance with no understanding of temporal relevancy. Regardless of approach, the outcome is the same: experts from future layers get evicted prematurely. We show this in Figure 3-top where the immediate next layer experts are most likely to be evicted. This observation motivates our Least-Stale approach, which exploits the spatial-temporal structure of expert access to prevent premature evictions.

#### 4.2 THE LEAST-STAILE POLICY

Eviction policies determine which cached experts to evict when capacity is exceeded. In many instances, eviction policies can cause **collision misses**, a type of cache miss where experts are evicted and then immediately needed again in the same forward pass. To address collision misses, Least-Stale (LS) combines both temporal factors (access time) and spatial awareness (layer positioning) to minimize collision misses.

*Staleness and position.* An expert is considered *stale* if it was accessed in a previous forward cycle; otherwise, it is *current*, meaning it is used in the current pass or selected by prefetching. We describe an expert as *left* if it appears earlier than the current execution point in layer-expert order.

Figure 3-bottom illustrates how staleness and position jointly determine eviction priority. The cache is partitioned into a *stale queue* containing experts from previous forward passes and a *current queue* containing experts from the ongoing pass.

Eviction prioritizes stale experts, as recently used experts are more likely to be reused by upcoming tokens or refreshed by prefetching. Experts in the current queue are evicted only to prevent out-of-memory errors. Within each queue, experts follow FIFO ordering based on layer position, favoring spatial and temporal locality.

In practice, both queues are implemented as priority heaps, requiring  $\mathcal{O}(N)$  space,  $\mathcal{O}(\log N)$  insertion, and  $\mathcal{O}(1)$  eviction, where  $N$  is the cache capacity.

### 5 EVALUATION AND RESULTS

Having motivated and described Least-Stale, we now present comprehensive evaluation results that validate our design choices and characterize performance across the full policy space.

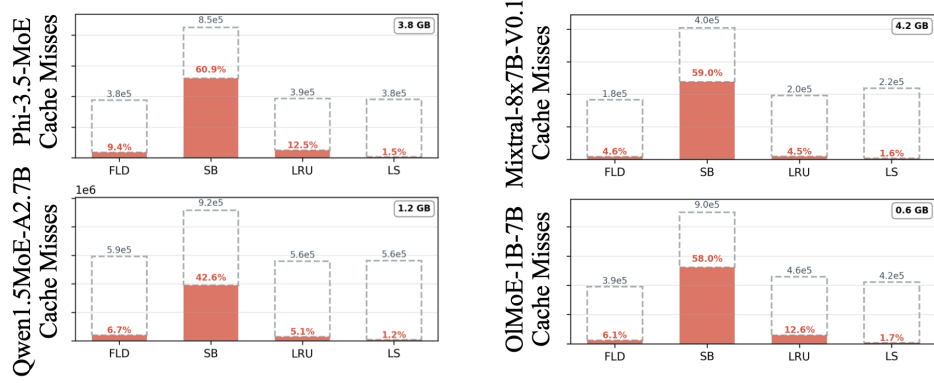


Figure 4: Collision miss across four eviction policies (FLD, SB, LRU, LS) at 5% capacity threshold showing total misses (gray) and collision misses (red). Refer to Section 3.1.3 for abbreviations.

### 5.1 EXPERIMENTAL SETUP

To effectively traverse the huge policy decision-space across various models and capacity constraints, we adopt a staged evaluation strategy: we first isolate optimal configuration from performance-neutral policies, i.e. eviction policies and prefetch policies, then evaluate quality-speed trade-offs using optimal configurations from early stages. This approach reduces the search space while ensuring fair comparison.

- **Models.** We evaluate across four MoE architectures: OLMoE-1B-7B Muennighoff et al. (2025), Mixtral-8x7B Jiang et al. (2024), Qwen1.5-MoE-A2.7B Team (2024), and Phi-3.5-MoE Abdin et al. (2024). Our selection covers both wide MoE with many small experts and deep MoE with a few yet large experts, enabling us to characterize how policy effectiveness varies with a model structure.

- **Tasks.** To measure realistic caching behaviour (with both prefill and decoding stages) we test tasks using autoregressive generation. These tasks are: GSM8K Cobbe et al. (2021) (mathematical reasoning), TruthfulQA Lin et al. (2021) (factuality), and NaturalQuestions Kwiatkowski et al. (2019) (question answering). Performance is measured via exact match for GSM8K and GPT-5 as judge for QA tasks. During policy benchmarking, we sampled only 100 examples for each task.

- **Hardware emulation.** We evaluate all models on a single A100 GPU with 80GB VRAM. We software-limit our GPU cache to 1%, 5%, or 25% of full model size, which amounts for example, to at most 21GB and as few as 0.8GB for Mixtral. We opted to use device native bandwidth, which on our device measures on average 5 GB/s. We quantize experts into int8 and int4 using bitsandbytes Dettmers et al. (2022) and optionally int2 using HQQ Badri & Shaji (2023) to evaluate cache miss policy comprehensively. For detailed model specifications please see in Appendix E.

### 5.2 LEAST-STALE EFFECTIVENESS

We first validate the Least-Stale policy by analyzing its impact on collision misses and inference speed across different models and capacity levels. We can see how baseline eviction policies at various capacities scored higher collision rates in Figure 4 which presents a 5% cache capacity. More extensive results can be found in Figure 8 of the Appendix.

For 5% cache capacity in particular, Least-Stale achieves 1.6-1.9% collision rates, while LRU reaches 4.5-12.6% and SB suffers catastrophic 42.6-60.9% rates. SB performs worst because historical gate scores capture neither spatial position nor temporal relevance; a highly-scored expert from layer 2 gets retained over a lower-scored expert from layer 10, even though layer 10 is spatially next. FLD, unlike LRU, tracks only expert distance from current layer and performed only slightly better than LRU.

- **Spatial-Temporal awareness.** Least-Stale determines eviction priority by forward-pass cycles. Experts accessed in the *current* pass are protected to prevent premature eviction of prefetched future experts and to increase the odds of reuse in the next forward cycle for left experts. Experts from *previous* passes are marked as stale and prioritized for eviction, unless they are needed again via prefetching.

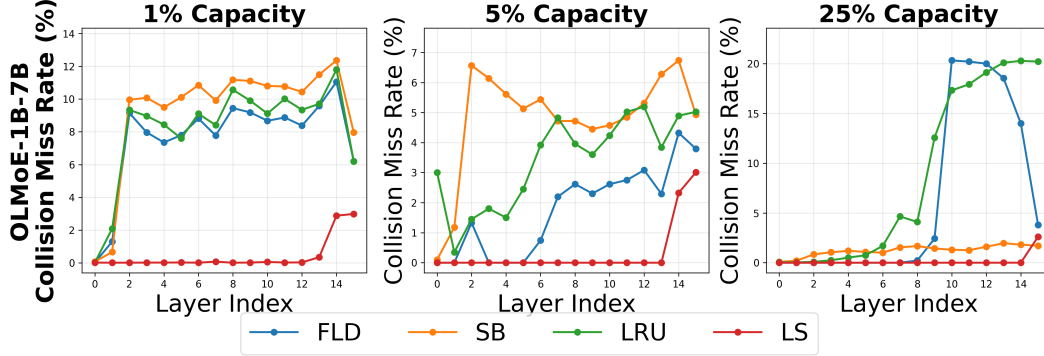


Figure 5: Per-layer collision miss rates for OLMoE across 16 layers. Least-Stale (red) maintains near-zero collisions across layers.

In Figure 5, we illustrate the per-layer collision misses. Baseline policies tend to accumulate throughout the forward pass as they evict soon-needed experts. At 1% capacity, LRU and Score-Based reach 8-12% collision rates in deeper layers. Least-Stale maintains near-zero collisions across all layers by only evicting stale left-layer experts.

• **Speed Impact.** The collision reduction translates directly to TTFT gains. At 5% capacity, Least-Stale reduces collision misses by 2.6-8.6 $\times$  versus LRU and 58-91 $\times$  versus Score-Based. When used as a drop-in replacement for existing approaches’ eviction policies (keeping their prefetch and routing unchanged), Least-Stale achieves 10.7-34.7% TTFT improvements on OLMoE (Figure 1), with consistent gains of 14.8-44.1% across diverse baseline systems (Table 2).

### 5.3 PREFETCH STRATEGY EFFECTIVENESS

#### Key Insight: Prediction accuracy $\neq$ cache performance

Score-based prefetching outperforms top-k despite lower overall prediction accuracy. Giving each layer the ability to adapt the number of experts to prefetch based on gating scores proves more effective at utilizing available bandwidth than a fix threshold.

• **Understanding the mismatch.** Using Least-Stale as our baseline eviction policy, we compare two prefetch strategies: top-k prefetching with different multiplier rates and score-based percentile filtering. What we discover is that higher prediction accuracy does not guarantee better cache performance in term of speed.

In Figure 6 we show that top-k with a 1.0 factor achieves high precision and recall across all models. Yet the right column shows that score-based prefetching typically delivers higher hit rates with lower synchronous overhead (time spent waiting on cache misses) for the majority of the model we tested, with the exception of Mixtral.

• **Effective Bandwidth Utilization.** We find that the score-based approach more effectively adapts the number of experts to prefetch on a per-layer basis. Simply scaling up a fixed number of experts does not match the same performance. We conclude that adaptive prefetching more effectively allocates existing bandwidth to high-reward layers while reducing the number of experts prefetched for high-risk layers, thereby outperforming fixed top-k strategies in most cases.

• **The Exception.** We observe Mixtral as the exception. It is a deep MoE with a small total number of experts, but each expert is large ( $\sim 336\text{MB}$ ). For these types of MoEs, the simpler top-k strategy performs competitively. This occurs because the size of the experts is too large for the available bandwidth, which limits the benefit of a dynamic prefetch system. This suggests that deep MoEs with large expert sizes may be better suited to simpler top-k prefetching.



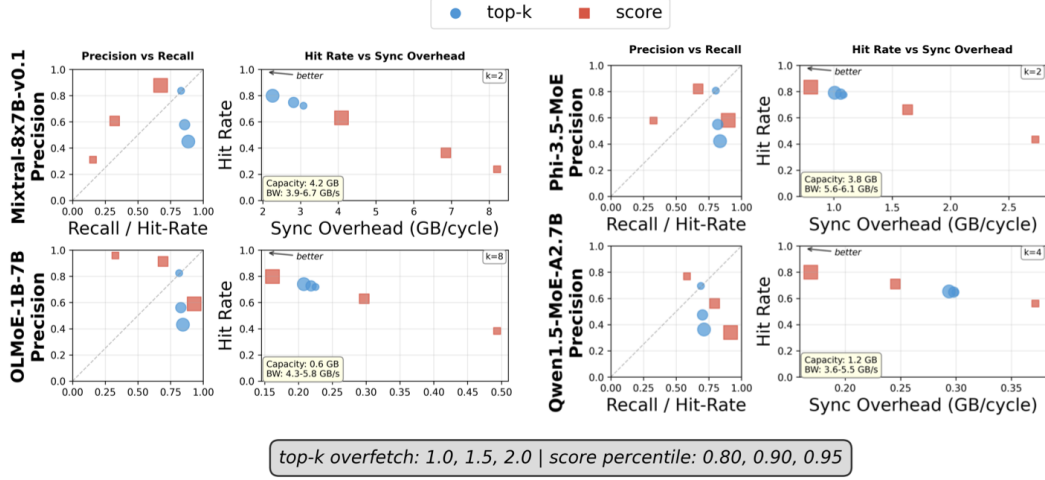


Figure 6: Prefetch strategy comparison at 5% capacity. Left: Precision vs recall shows top-k 1.0 $\times$  achieves highest prediction accuracy. Right: Hit rate vs synchronous overhead reveals score-based 80th percentile achieves better cache performance for most models through dynamic adaptation.

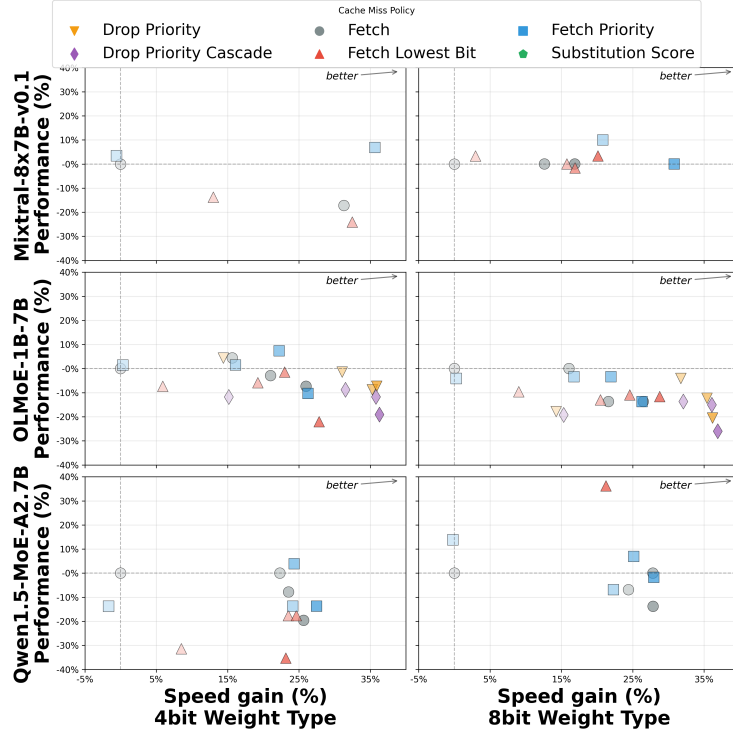


Figure 7: Quality-memory trade-offs for six miss handling policies at 5% capacity. Performance and speed shown relative to Fetch policy baseline. We report the average accuracy among the tasks under study as performance. Points above the zero horizontal line indicate performance improvements compared to the baseline; points to the right of the zero vertical line indicate speed gains. Color intensity indicates cache-aware routing strength ( $\lambda$ ). Substitution policies fall outside the chart range.

Table 1: Baseline comparison across all models at 5% cache capacity (4bit quantization, 5 GB/s bandwidth, with QKV Cache). TTFT in milliseconds. SpecMD is capable of exploring multiple configurations of the state-of-the-art policies, and one of them is our chosen configuration (Config 5) from Figure 7, balancing between performance and overall speed.

Approach		OLMoE-1B-7B					Mixtral-8x7B				
		GSM8K	Truth.	NQ	TTFT	Token/s	GSM8K	Truth.	NQ	TTFT	Token/s
Baseline (8bit)		0.704	0.540	0.361	3166	1.5	0.569	0.345	0.510	10296	0.4
Config#1	Spec <sup>+</sup> , LRU <sup>◦</sup> , Fetch <sup>□</sup> , Stnd <sup>§</sup>	0.677	0.555	<b>0.343</b>	3077	1.91	<b>0.566</b>	0.343	<b>0.518</b>	4238	0.98
Config#2	Score <sup>♣</sup> , ScoreE <sup>♣</sup> , Subs <sup>◇</sup> , Stnd	0.073	0.493	0.159	2257	1.22	0.037	0.251	0.273	<b>1620</b>	0.75
Config#3	Spec, LHU <sup>▽</sup> , Cascade <sup>△</sup> , Stnd	0.606	0.548	0.291	<b>1590</b>	1.34	0.310	0.326	0.347	2322	0.71
Config#4	No Pref, LRU, Fetch, Mod <sup>ψ</sup>	0.664	0.558	0.332	2665	2.33	<b>0.566</b>	0.326	0.517	4333	1.07
Config#5	Score, LS, Fetch, Stnd	<b>0.691</b>	<b>0.559</b>	0.331	2220	1.69	0.528	<b>0.356</b>	0.471	3044	1.42

Approach		Phi-3.5-MoE					Qwen1.5-MoE-A2.7B				
		GSM8K	Truth.	NQ	TTFT	Token/s	GSM8K	Truth.	NQ	TTFT	Token/s
Baseline (8bit)		0.697	0.774	0.594	7902	0.8	0.471	0.595	0.460	5728	2.3
Config#1	Spec <sup>+</sup> , LRU <sup>◦</sup> , Fetch <sup>□</sup> , Stnd <sup>§</sup>	<b>0.705</b>	0.775	0.589	4452	1.69	<b>0.506</b>	<b>0.614</b>	0.473	3724	2.89
Config#2	Score <sup>♣</sup> , ScoreE <sup>♣</sup> , Subs <sup>◇</sup> , Stnd <sup>§</sup>	0.027	0.514	0.269	<b>1906</b>	0.86	0.008	0.480	0.408	2739	1.11
Config#3	Spec, LHU <sup>▽</sup> , Cascade <sup>△</sup> , Stnd	0.604	0.743	0.563	3459	1.05	0.217	0.608	0.462	2492	2.19
Config#4	No Pref, LRU <sup>◦</sup> , Fetch <sup>□</sup> , Mod <sup>ψ</sup>	<b>0.705</b>	0.777	<b>0.591</b>	5025	1.72	0.296	0.611	0.465	3489	3.50
Config#5	Score, LS, Fetch, Stnd	0.684	<b>0.780</b>	0.587	3863	1.11	0.497	0.607	<b>0.474</b>	2893	2.11

<sup>◦</sup> Speculative Prefetch; <sup>□</sup> LRU Eviction Eliseev & Mazur (2023); <sup>§</sup> Fetch immediately for miss; <sup>§</sup> Standard routing; <sup>♣</sup> Score-based Prefetch Zhu et al. (2025); <sup>♣</sup> Score-based Eviction Zhu et al. (2025); <sup>◇</sup> Substitution with cached experts for miss Zhu et al. (2025); <sup>▽</sup> Least High-precision Eviction Tang et al. (2024); <sup>△</sup> Fetch with precision cascade; <sup>ψ</sup> Routing Modification Skliar et al. (2025)

#### 5.4 QUALITY-SPEED TRADE-OFFS

• **Trading Quality for Speed.** With optimal eviction and prefetch policies established, we now examine how different miss handling and routing strategies perform. By modifying these two policies, we can achieve a quality/speed tradeoff, shown in Figure 7.

when studying miss-handling policies, we find that Fetch Priority (blue squares) consistently occupies the Pareto frontier across models. This approach stores multiple precision levels and selectively degrades low-importance experts (bottom 60th percentile) during cache misses. For Mixtral at 8-bit, we observe 10%-12% performance improvements with 50-60% speed gains. For OLMoE at 4-bit, Fetch Priority maintains near-baseline performance (0-5% loss) while achieving 20%-30% speedup.

• **The Trade-offs.** Drop Priority, as expected, provides the highest speed improvements at the cost of performance degradation. Qwen at 4-bit suffers -25% to -30% performance loss, while OLMoE, a wide MoE, shows only -5% to -15% degradation. This effectively demonstrates that more active experts can dilute the impact of dropping any single expert. We find that Fetch Lowest Bit and Drop Priority Cascade show inconsistent results across models, offering no clear advantage over simpler strategies. On the other hand, Substitution Score fails across all configurations, suggesting that replacing missing experts with functionally similar alternatives is unreliable, regardless of capacity level.

• **Cache-Aware Routing Impact.** We find that cache-aware routing is credited for 10%-20% of the achieved speed improvements in some cases (compare different intensities of the blue squares on each figure). The effectiveness depends critically on model architecture: OLMoE can tolerate higher routing bias ( $\lambda$ ) than Mixtral, since wider expert distributions can redistribute the impact of routing perturbations more effectively. Our experiments show that  $\lambda$  values between 0.1-0.5 work best, balancing speed gains with quality preservation; too conservative values provide minimal benefit, too aggressive values degrade quality unacceptably.

#### 5.5 SPECMD MIX-POLICY EVALUATION

SpecMD’s enables an easy way to mix and run ad-hoc policies across all four dimensions (eviction, prefetch, miss handling, routing) on the same testing platform. Through it, we can identify policies with favorable characteristic via policy evaluation. In previous sub-sections, we discussed at length

policy interactions and the various tradeoffs; we now selected five set of mix policies to compare running full task evaluation. This we summarize in Table 1 using 5% capacity cache capacity and 4bit quantization.

- **Configuration selection methodology.** To ensure fair comparison and avoid cherry-picking, we first select a four configurations (#1-#4) obtained from our previous evaluation (Section 5.3) that characterized the pareto curve with enough policy diversity. We then select *a configuration* from the same study that scores high on both performance and speed for most models.

- **Results interpretation.** Configurations selected using SpecMD perform well in most cases, but this is not the point. The point of SpecMD is the experimental platform which enables flexible configuration and unbiased study of the policies and their interactions. Using SpecMD, researchers can select optimal policy combinations based on their specific hardware constraints and quality requirements and share their findings more effectively.

## 6 CONCLUSION

This work presented SpecMD, a unified framework for evaluating MoE expert caching policies. Through GPU-based hardware emulation and minimal PyTorch integration, we enabled controlled comparison of caching strategies without specialized hardware requirements.

Our study challenges traditional caching assumptions for MoE architectures. We demonstrated that expert access follows deterministic sequential layer patterns rather than temporal locality, rendering conventional LRU/LFU approaches ineffective. The proposed Least-Stale eviction policy exploits this structure, achieving 3-9× collision miss reduction at practical capacity levels, translating as an average 17% improvement in TTFT, across different models and methods. Counterintuitively, we revealed that prediction accuracy does not guarantee cache performance—score-based prefetching delivers superior hit rates through dynamic bandwidth adaptation despite lower precision metrics.

## IMPACT STATEMENT

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## REFERENCES

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone, 2024. URL <https://arxiv.org/abs/2404.14219>.

- Hicham Badri and Appu Shaji. Half-quadratic quantization of large machine learning models, November 2023. URL [https://mobiusml.github.io/hqq\\_blog/](https://mobiusml.github.io/hqq_blog/).
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Hao Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang Yuan, Junjie Qiu, Junxiao Song, Kai Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhua Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Size Zheng, T. Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Liu, Xin Xie, Xingkai Yu, Xinnan Song, Xinyi Zhou, Xinyu Yang, Xuan Lu, Xuecheng Su, Y. Wu, Y. K. Li, Y. X. Wei, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Zheng, Yichao Zhang, Yiliang Xiong, Yilong Zhao, Ying He, Ying Tang, Yishi Piao, Yixin Dong, Yixuan Tan, Yiyuan Liu, Yongji Wang, Yongqiang Guo, Yuchen Zhu, Yudian Wang, Yuheng Zou, Yukun Zha, Yunxian Ma, Yuting Yan, Yuxiang You, Yuxuan Liu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhewen Hao, Zhihong Shao, Zhiyuan Wen, Zhipeng Xu, Zhongyu Zhang, Zhuoshu Li, Zihan Wang, Zihui Gu, Zilin Li, and Ziwei Xie. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model, 2024. URL <https://arxiv.org/abs/2405.04434>.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- Artyom Eliseev and Denis Mazur. Fast inference of mixture-of-experts language models with offloading. *ArXiv*, abs/2312.17238, 2023. URL <https://api.semanticscholar.org/CorpusID:266573098>.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L  lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th  ophile Gerv  t, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mixtral of experts, 2024. URL <https://arxiv.org/abs/2401.04088>.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*, 2019.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods, 2021.
- Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, Yuling Gu, Shane Arora, Akshita Bhagia, Dustin Schwenk, David Wadden, Alexander Wettig, Binyuan Hui, Tim Dettmers, Douwe Kiela, Ali Farhadi, Noah A. Smith, Pang Wei Koh, Amanpreet Singh, and Hannaneh Hajishirzi. Olmoe: Open mixture-of-experts language models, 2025. URL <https://arxiv.org/abs/2409.02060>.

Andrii Skliar, Ties van Rozendaal, Romain Lepert, Todor Boinovski, Mart van Baalen, Markus Nagel, Paul Whatmough, and Babak Ehteshami Bejnordi. Mixture of cache-conditional experts for efficient mobile device inference, 2025. URL <https://arxiv.org/abs/2412.00099>.

Peng Tang, Jiacheng Liu, Xiaofeng Hou, Yifei Pu, Jing Wang, Pheng-Ann Heng, Chao Li, and Minyi Guo. Hobbit: A mixed precision expert offloading system for fast moe inference. *ArXiv*, abs/2411.01433, 2024. URL <https://api.semanticscholar.org/CorpusID:273812295>.

Qwen Team. Qwen1.5-moe: Matching 7b model performance with 1/3 activated parameters”, February 2024. URL <https://qwenlm.github.io/blog/qwen-moe/>.

Guoying Zhu, Meng Li, Haipeng Dai, Xuechen Liu, Weijun Wang, Keran Li, Jun xiao, Ligeng Chen, and Wei Wang. Enabling moe on the edge via importance-driven expert scheduling. *ArXiv*, abs/2508.18983, 2025. URL <https://api.semanticscholar.org/CorpusID:280869816>.

## A IMPLEMENTATION DETAILS

### A.1 GATE HOOK ARCHITECTURE

Our framework intercepts MoE gate layer outputs through PyTorch forward hooks. When a gate layer executes, it outputs a 4-tuple: selected expert indices, routing weights for those experts, raw router logits, and normalized router weights across all experts. We wrap gate outputs to ensure consistent formatting across different MoE architectures (Mixtral, OLMoE, Qwen), providing a unified interface for downstream components. This standardization eliminates model-specific handling in policy implementations.

The hook manager registers forward hooks on each gate layer to intercept these outputs. Algorithm 1 (main text) outlines the core hook logic. Upon receiving gate outputs, the routing policy can modify logits to prefer GPU-resident experts. When routing modifies expert selections, we apply dual-logit semantics: modified logits determine which experts to select, but original logits determine their contribution weights. This preserves model output integrity—even if cache-aware routing changes which experts execute, their contributions are weighted according to the original router decisions.

### A.2 CACHE MANAGER DETAILS

For each selected expert during model execution, the cache manager checks GPU cache availability. Cache hits proceed immediately. Cache misses trigger the miss handling policy: fetch blocks until the expert loads from CPU, drop skips low-rank experts, or substitution replaces missing experts with cached alternatives. The cache manager maintains experts across CPU (pinned memory) and GPU (active cache) storage, supporting mixed-precision configurations where experts exist in multiple quantization levels (fp16, int8, int4, int2).

### A.3 PREFETCH MANAGER AND WATCHDOG THREAD

Concurrently, the prefetch manager anticipates future expert needs. The hook manager uses either current gate outputs or future gate functions on current gate logits to predict which experts subsequent layers will require, submitting prefetch requests to an asynchronous queue. A watchdog thread monitors this queue and GPU capacity, issuing non-blocking CUDA transfers when space allows. Prefetched experts populate the GPU cache before they are requested, reducing miss frequency. When capacity pressure exceeds limits, the eviction policy identifies victims—our Least-Stale policy prioritizes removing experts from already-processed layers while protecting experts needed for upcoming layers.

### A.4 HARDWARE EMULATION

Hardware constraints are emulated through software limits rather than requiring physical low-bandwidth devices. Bandwidth throttling delays expert transfers between CPU and GPU to simulate slower interconnects. Capacity limits restrict the available GPU cache below physical VRAM, forcing the system to make eviction and prefetch trade-offs. This emulation enables systematic exploration: evaluating a policy under 5 GB/s bandwidth and 25% capacity requires only configuration changes, not specialized hardware.

## B BASELINE EVICTION POLICIES

We implement five baseline eviction strategies for comparison:

**LRU (Least Recently Used)** evicts the least recently accessed expert. However, MoE expert access follows deterministic sequential layer patterns rather than temporal patterns, making LRU ineffective for expert caching.

**LFU (Least Frequently Used)** evicts the expert with the lowest access count over the entire inference history. This policy prioritizes retaining frequently-accessed experts across all layers. However, it fails to account for layer-wise access patterns where experts from completed layers will not be needed until the next forward pass.

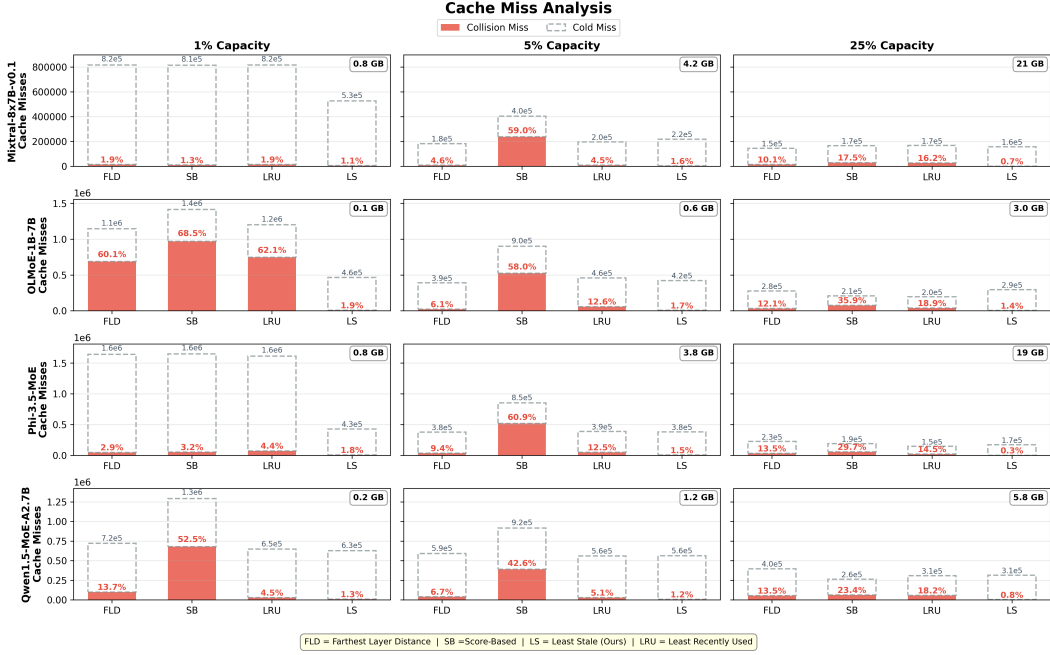


Figure 8: Complete collision miss analysis across four eviction policies (FLD, SB, LRU, LS) at three capacity levels (1%, 5%, 25%). Least-Stale (LS) achieves dramatically lower collision miss rates (red bars) compared to baseline policies across all models and capacity constraints. The pattern is consistent: LS achieves 10-34 $\times$  lower collision rates than baseline policies.

**LHU (Least High-precision Used)** Tang et al. (2024) tracks high-precision expert usage frequency, prioritizing retention of frequently-accessed high-precision experts since cache misses for these incur significantly higher latency. This policy must be combined with fetch cascade precision policies (e.g., int8 $\rightarrow$ int4 $\rightarrow$ int2) that allow graceful degradation to lower precision on cache misses.

**FLD (Farthest Layer Distance)** Tang et al. (2024) exploits layer-wise structure by retaining experts from nearer layers, which are more likely to be needed soon. This policy computes the distance between each cached expert’s layer and the current inference layer, evicting experts from the farthest layers first.

**Score-based** Zhu et al. (2025) evicts experts with the lowest historical gate scores, using router activation values as importance signals. This policy maintains a running history of gate scores for each expert and evicts those with minimal contribution to model outputs. However, historical scores do not capture position in the forward pass, causing premature eviction of soon-needed experts.

## C COMPLETE EVICTION POLICY RESULTS

Figure 8 presents complete collision miss rate results referenced in Section 5.1, showing consistent performance across all capacity levels.

Figure 9 shows per-layer collision patterns for all models, revealing the mechanism behind Least-Stale’s effectiveness across different architectures.

## D LEAST-STATE EVICTION AS DROP-IN IMPROVEMENT

Beyond our full system, we evaluate the Least-Stale eviction policy as a standalone contribution that can improve existing approaches. As visualized in Figure 1 for OLMoE-1B-7B, Least-Stale provides consistent speedups across diverse baseline systems. Table 2 shows comprehensive TTFT improvements across all four models when we replace each baseline’s original eviction policy with Least-Stale while keeping their other components (prefetching, routing) unchanged.

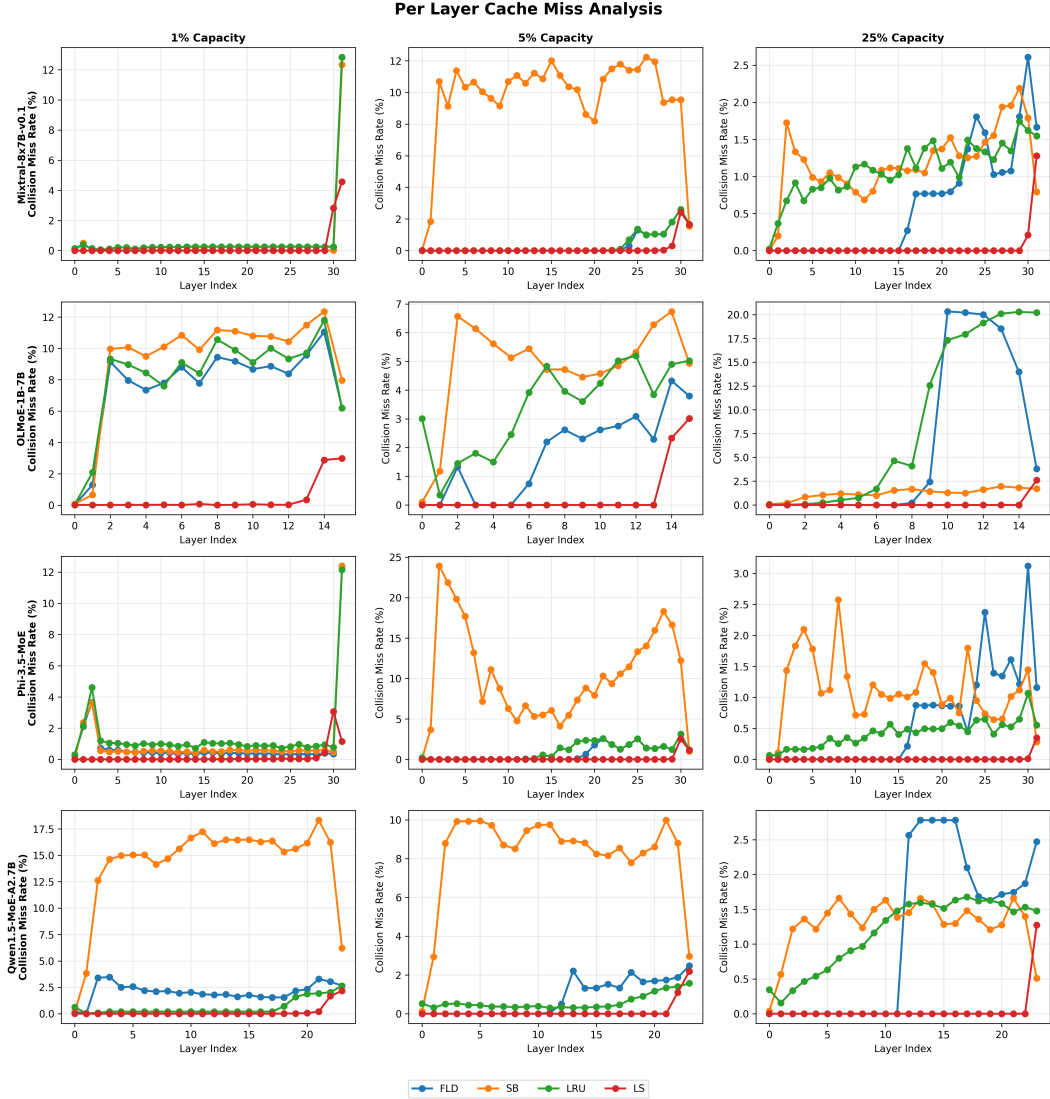


Figure 9: Per-layer collision miss analysis across all four models at three capacity levels (1%, 5%, 25%). Least-Stale (red) maintains near-zero collisions by protecting experts in the current forward pass, while baseline policies show accumulating collisions throughout inference. The mechanism reveals why Least-Stale achieves dramatic improvements: it only evicts stale left-layer experts while baseline policies prematurely evict soon-needed experts.

Least-Stale eviction provides consistent TTFT improvements across most configurations (10-44% reduction), demonstrating its value as a portable component. However, SpecMD’s advantage stems from the synergistic combination of eviction, prefetching, and cache-aware routing policies.

## E MODEL SPECIFICATIONS

Table 3 provides detailed specifications for all models evaluated in this work.

These specifications highlight the diversity of architectural choices: OLMoE uses many small experts with high top-k (fine-grained sparsity), while Mixtral uses few large experts with low top-k (coarse-grained sparsity). This diversity enables us to characterize policy effectiveness across different granularity regimes.



Model	Approach	Original TTFT (ms)	w/ Least-Stale (ms)	Improvement (%)
<i>OLMoE-1B-7B</i>				
	Eliseev '23	3076.8	2748.1	10.7%
	Zhu et al. '25	2257.2	1473.1	34.7%
	HOBBIT '24	1590.1	1309.2	17.7%
	MixtureCache '25	2664.6	2297.5	13.8%
<i>Mixtral-8x7B</i>				
	Eliseev '23	4238.0	3485.2	17.8%
	Zhu et al. '25	1620.3	1460.8	9.8%
	HOBBIT '24	2321.7	1725.2	25.7%
	MixtureCache '25	4332.6	3573.2	17.5%
<i>Phi-3.5-MoE</i>				
	Eliseev '23	4451.7	3791.1	14.8%
	Zhu et al. '25	1905.6	1987.6	-4.3%
	HOBBIT '24	3458.6	1933.6	44.1%
	MixtureCache '25	5025.3	3791.1	24.6%
<i>Qwen1.5-MoE-A2.7B</i>				
	Eliseev '23	3724.0	3433.4	7.8%
	Zhu et al. '25	2739.2	2620.9	4.3%
	HOBBIT '24	2492.0	2033.7	18.4%
	MixtureCache '25	3488.6	2685.4	23.0%

Table 2: Least-Stale eviction policy as drop-in improvement for existing approaches. We replace each baseline’s original eviction policy with Least-Stale while keeping their prefetching and routing strategies unchanged. Positive percentages indicate TTFT reduction (lower is better). This demonstrates Least-Stale’s generalizability across different system architectures.

Table 3: Detailed specifications for evaluated MoE architectures. All sizes in GB unless specified. Activated size represents memory required per forward pass.

Model	Layers	Experts per Layer	Top-k	Expert Size	Total Size	Expert Params	Non-Expert	Activated
OLMoE-1B-7B	16	64	8	12 MB	12.89 GB	12.0 GB	0.89 GB	2.39 GB
Mixtral-8x7B	32	8	2	336 MB	86.99 GB	84.0 GB	2.99 GB	23.99 GB
Qwen1.5-MoE-A2.7B	24	60	4	16.5 MB	26.67 GB	23.2 GB	3.46 GB	5.01 GB
Phi-3.5-MoE	32	16	2	150 MB	77.99 GB	75.0 GB	2.99 GB	12.37 GB