# Stateful Cache-Aware Routing for Distributed LLM Inference:
# A Push-Based Consistency Approach

**Team:** Aman Pandey, Sanjeev Kumar, Aryan Khanolkar, Guna Avula
*Purdue University*

November 2025

### Abstract

The efficiency of Large Language Model (LLM) serving is fundamentally constrained by the memory-bound "prefill" phase. While execution engines like vLLM optimize intra-node memory via PagedAttention, distributed deployments typically rely on stateless load balancing. This disconnect leads to "Cache Blindness," where sequential requests sharing identical system prompts or retrieval contexts are scattered across replicas, forcing redundant re-computation.

We present a **Stateful Cache-Aware Router** that bridges the gap between local engine optimizations and cluster-wide efficiency. Unlike existing solutions that rely on high-latency disaggregated storage (LMCache) or heavy orchestration (llm-d), our system introduces a novel **Push-Based Eviction Callback** protocol. By instrumenting vLLM workers to actively report cache lifecycle events and implementing an **Anti-Entropy State Reconciliation** mechanism, our architecture guarantees eventual consistency with minimal overhead. Benchmarks demonstrate that this approach effectively eliminates redundant prefills in Retrieval-Augmented Generation (RAG) workloads.

## 1  Introduction

The computational cost of Generative AI is dominated by the attention mechanism, which scales quadratically $O(L^2)$ with sequence length $L$. In modern RAG workloads, input contexts frequently exceed 10,000 tokens, making the "prefill" phase the primary latency bottleneck. **Automatic Prefix Caching (APC)** has emerged as a standard optimization within inference engines, allowing Key-Value (KV) tensors to be persisted and reused across requests.

However, APC is inherently local to a single GPU. In a distributed cluster of $N$ replicas, a standard Round-Robin load balancer offers a cache hit probability of only $1/N$ for a repeated prefix. This project addresses the **"Distributed Locality Gap"**—the failure of stateless routing to leverage stateful engine optimizations.

We propose a middleware solution that adheres to the principle of **"Compute-to-Data"** routing. Rather than moving heavy KV states across the network, we route lightweight requests to the worker that already holds the state.

## 2  Design Rationale

The inception of this system stems from a critical observation: **There is a disconnect between the "Compute Layer" (vLLM) and the "Routing Layer" (Load Balancer).**

1

## 2.1   The Logical Progression

1. **The Problem**: vLLM is excellent at *local* optimization. If a request hits a worker that has the prefix cached, it's fast. But standard load balancers (Round-Robin, Least-Connection) are "stateless"—they don't know which worker has which prefix.

2. **The Consequence**: This leads to "Cache Blindness." A request with a heavy system prompt might be routed to Worker B, even though Worker A just served it and has the cache hot. Worker B has to recompute the KV cache (expensive prefill), wasting GPU cycles and increasing latency.

3. **Existing Solutions & Their Flaws**:

   - **Shared Storage (LMcache)**: "Let's move the cache to a shared disk/Redis." → **Flaw**: Moving gigabytes of KV cache over the network is slow. It introduces latency that often negates the benefit for short-to-medium prompts.
   - **Heavy Orchestration (llm-d)**: "Let's build a complex Kubernetes operator." → **Flaw**: High barrier to entry. Tightly coupled to K8s. Hard to debug.

4. **Our "Aha!" Moment**: We don't need to move the data (LMcache), and we don't need a heavy platform (llm-d). We just need a **lightweight signal** to tell the router where the data *already is*.

## 2.2   Why "Push-Based" Eviction?

Most systems use "Pull" (polling) or "Estimation".

- **Pull**: The router asks workers "What do you have?" → Too much traffic, always outdated.

- **Estimation**: The router guesses "I sent it to Worker A, so Worker A must have it." → **Flaw**: What if Worker A ran out of memory and evicted it 10ms ago? The router would send a "False Hit."

**Our Solution**: **Push-Based Eviction Callbacks**.

- The worker is the *source of truth*.

- When vLLM evicts a block, it *immediately* notifies the router.

- This ensures the router's map is "Eventually Consistent" with very low latency.

# 3   System Architecture

Our system employs a three-tier architecture designed for modularity and fault tolerance.

## 3.1   Tier 1: The Tokenizer-Aware Router

The router is the control plane. It does not perform inference but makes scheduling decisions based on a `GlobalCacheMap`. Crucially, it must replicate the hashing logic of the inference engine exactly.

- **Tokenization Alignment:** Hashing raw strings is unreliable (e.g., whitespace sensitivity). The router loads the specific model tokenizer (e.g., Mistral-7B) to convert prompts to Token IDs before hashing.

- **Prefix-Aware Hashing:** To support RAG/Chat, where a system prompt is constant but the user query varies, the router calculates the hash based on an `effective_prefix_len`.

## 3.2  Tier 2: The Worker Pool (vLLM + Sidecar)

We modify standard vLLM instances by injecting a surgical "sidecar" thread called the `EvictionReporter`. This component hooks into the `BlockManager.free_block()` method. Unlike polling systems, this makes the worker the **Source of Truth**.

## 3.3  Tier 3: The Consistency Protocol

To maintain $V \approx S$, we implement a hybrid consistency protocol:

1. **Fast Path (Eviction Reporting):** Real-time notification of cache removals. To prevent Distributed Denial of Service (DDOS) on the router during mass evictions, these reports are batched and flushed every 100ms.

2. **Slow Path (Anti-Entropy Sync):** Periodic background reconciliation (every 5s) where workers send a full snapshot of their active state. This heals divergence caused by network partitions.

# 4  Deep Dive: Implementation Details

## 4.1  Router: Tokenizer-Aware Hashing

The core correctness requirement is that $\text{Hash}_{\text{Router}}(P) \equiv \text{Hash}_{\text{Worker}}(P)$. If the router and worker disagree on the hash of a prefix, requests will be routed incorrectly (0% Hit Rate). We solve this by ensuring the router performs the exact same tokenization as the worker before hashing.

```
class TokenizerUtils:
    def compute_prefix_hash(self, text: str, prefix_len: int = None) ->
    str:
        token_ids = self.tokenize(text)
        # Use effective length to ignore variable user queries
        if prefix_len:
            token_ids = token_ids[:prefix_len]

        # Hash the tuple of integers for stability.
        # This matches vLLM's internal block hashing strategy.
        return hashlib.sha256(str(tuple(token_ids)).encode()).hexdigest
    ()
```

Listing 1: Robust Prefix Hashing Strategy

## 4.2  Router: Scalable Reverse Index

A naive implementation of the 'sync' operation would require iterating over the entire global map to remove a worker's stale entries. As the number of cached prefixes grows ($10^5+$), this becomes an $O(\text{TotalHashes})$ operation.

We implemented a **Reverse Index** optimization:

- **Forward Map:** `PrefixHash` $\rightarrow$ `Set[WorkerID]`

- **Reverse Map:** `WorkerID` $\rightarrow$ `Set[PrefixHash]`

This reduces the synchronization complexity to $O(\text{WorkerHashes})$, decoupling the router's CPU load from the total cluster size.

## 4.3   Worker: Asynchronous Batching

Hooking into the critical path of memory management requires extreme care. Blocking the GPU engine to send an HTTP request would degrade inference performance. We use a thread-safe deque for atomic hand-off between the GPU thread and the IO thread.

```python
# vLLM Hook (Pseudo-code)
def free_block(self, block):
    # Atomic push to memory queue (Microsecond latency)
    if block.is_cached_prefix:
        self.eviction_queue.append(block.hash)

# Background Reporter Thread (Actual Implementation)
async def _report_loop(self):
    async with aiohttp.ClientSession() as session:
        while self._running:
            await asyncio.sleep(self.report_interval)
            if not self.eviction_queue: continue

            # Drain the queue
            batch = []
            while self.eviction_queue:
                batch.append(self.eviction_queue.popleft())

            if not batch: continue

            try:
                payload = {"worker_id": self.worker_id, "evicted_hashes": batch}
                async with session.post(f"{self.router_url}/internal/eviction", json=payload) as resp:
                    if resp.status != 200:
                        logger.error(f"Failed to report evictions: {resp.status}")
            except Exception as e:
                logger.error(f"Error reporting evictions: {e}")
```

Listing 2: Non-Blocking Eviction Reporting (EvictionReporter)

## 4.4   Anti-Entropy State Reconciliation

Distributed systems must withstand partial failures. If a "Fast Path" eviction report is dropped due to a network glitch, the router will hold a "Phantom Cache" entry—believing data exists when it does not.

To solve this, we implemented a self-healing "Slow Path":

1. Every 5 seconds, the worker scans its internal block table to identify all active prefix hashes.

2. It sends this full list to the router via `/internal/sync`.

3. The router performs a **Destructive Update**: it clears all known entries for that worker and repopulates them from the list.

This guarantees that any state divergence persists for at most 5 seconds.

```python
# Router: cache_map.py
```

```python
def sync_worker_state(self, worker_id: str, active_hashes: List[str]):
    with self._lock:
        # 1. Destructive Update: Clear old state for this worker using
   Reverse Index
        # Complexity: O(M) where M is number of hashes held by THIS
   worker.
        if worker_id in self._worker_to_hashes:
            current_hashes = list(self._worker_to_hashes[worker_id])
            for h in current_hashes:
                if h in self._map:
                    self._map[h].discard(worker_id)
                    if not self._map[h]:
                        del self._map[h]
            self._worker_to_hashes[worker_id].clear()

        # 2. Re-populate with authoritative current state
        for h in active_hashes:
            if h not in self._map: self._map[h] = set()
            self._map[h].add(worker_id)

            if worker_id not in self._worker_to_hashes:
                self._worker_to_hashes[worker_id] = set()
            self._worker_to_hashes[worker_id].add(h)
```

Listing 3: Router State Reconciliation

# 5 Routing Algorithm

The router executes a **Sticky-Least-Loaded** policy with **Speculative Updates**.

---
**Algorithm 1** Stateful Routing with Speculation
---
1: **procedure** ROUTEREQUEST(*req*)
2:     $L \leftarrow req.prefix\_len$ **if** $req.prefix\_len$ **else** 64                                    ▷ Smart Default
3:     $h \leftarrow \text{ComputeHash}(req.prompt, L)$
4:     $candidates \leftarrow \text{GlobalMap}.get(h)$
5:     **if** $candidates \neq \emptyset$ **then**
6:                                                                     ▷ **Cache Hit**: Stickiness ensures locality
7:         **return** $\text{argmin}_{w \in candidates}(\text{Load}(w))$
8:     **else**
9:                                                                     ▷ **Cache Miss**: Load Balancing
10:         $target \leftarrow \text{argmin}_{w \in W}(\text{Load}(w))$
11:         $\text{GlobalMap}.update(h, target)$                                ▷ **Speculative Update**
12:         **return** $target$
13:     **end if**
14: **end procedure**
---

**Rationale for Speculation:** When a request is routed to a worker on a miss, we *immediately* record that the worker has the prefix. This prevents the "Thundering Herd" problem where concurrent requests for the same new prefix are routed to different workers because the first worker hasn't reported the cache creation yet.

# 6  Evaluation and Validation

## 6.1  Benchmark Design

We developed an end-to-end benchmarking suite ('benchmark.py') to validate the system lifecycle across four states:

1. **Cold Start:** Verifies Speculative Update.

2. **Stickiness:** Verifies that a subsequent request routes to the same worker, even if another worker has lower load.

3. **Eviction:** Verifies the worker correctly reports memory pressure.

4. **Recovery:** Verifies the router updates the map and routes the next request to a new worker.

## 6.2  Comparative Analysis

Table 1 summarizes our architectural position relative to SOTA.

| Feature | vLLM (Stock) | LMCache | llm-d | Ours |
|---|---|---|---|---|
| Routing Strategy | Round-Robin | N/A | K8s Gateway | **Sticky-Hash** |
| Data Locality | None | Low | High | **High** |
| State Consistency | N/A | Strong | Eventual (Poll) | **Eventual (Push)** |
| Network Overhead | Low | High | Medium | **Very Low** |
| Infra Dependency | None | Redis/Disk | Kubernetes | **None (Python)** |

Table 1: Comparison of Distributed Inference Architectures

# 7  Discussion: Limitations and Future Work

## 7.1  Failure Modes Analysis

- **Router Crash:** Since the map is in-memory, a router restart clears the state. The **Anti-Entropy Sync** loop in the workers ensures the router rebuilds its map within 5 seconds of coming back online.

- **Network Partition:** If eviction reports are dropped, the router may have "Phantom Cache" entries (False Hits). The Sync loop heals this automatically.

## 7.2  Roadmap for Production

1. **Persistence:** Moving 'GlobalCacheMap' to Redis would allow the router to be stateless and horizontally scalable.

2. **Security:** Internal endpoints ('/internal/*') must be secured via mTLS or VPC isolation in a production environment.

3. **Model Generalization:** The current hashing schema assumes a single model. Production systems should include 'model_name' in the hash key.

# 8  Conclusion

This project moves beyond theoretical design to provide a robust, implementation-ready solution for distributed LLM inference. By solving the critical engineering challenges of **Tokenization Alignment**, **Network Backpressure** (via batched reporting), and **State Consistency** (via anti-entropy sync), we provide a scalable path to reducing TTFT in production RAG clusters. The system effectively bridges the gap between local engine optimizations and global cluster routing.