

Critical Evaluation of Auxiliary Cache Utility Data.pdf

by Turnitin Student

Submission date: 09-Jun-2025 02:57AM (UTC-0400)

Submission ID: 2695277929

File name: Critical_Evaluation_of_Auxiliary_Cache_Utility_Data.pdf (307.63K)

Word count: 6180

Character count: 38856

Critical Evaluation of Auxiliary Cache Utility Data Structures in Vector Databases and LLM Retrieval Pipelines

Mateus Yonathan
Software Developer & Independent Researcher

<https://www.linkedin.com/in/siyoyo/>

Abstract—This paper presents an analytical framework for evaluating auxiliary cache utility data structures in vector databases and local large language model (LLM) retrieval pipelines. While frequency-based caching policies such as TinyLFU are well-established in mature systems, the incremental value of a dedicated auxiliary cache layer that sits between application logic and existing caching mechanisms remains an open question. We develop a formal mathematical model for analyzing cache utility and present a structured assessment methodology examining architectural trade-offs. Our analysis encompasses system complexity, maintenance burden, and cache coherence challenges against potential performance benefits. We systematically compare architectural approaches across different deployment scenarios and query workloads, identifying conditions where auxiliary caches might provide meaningful benefits beyond existing mechanisms. Rather than making empirical claims that would require extensive benchmarking, we contribute a decision framework that system architects can apply to their specific contexts. We conclude by identifying key experimental metrics and validation approaches that would be necessary to quantitatively evaluate auxiliary cache implementations in production environments.

I. INTRODUCTION

As organizations increasingly adopt large language models (LLMs), the efficiency and reliability of retrieval mechanisms have become critical concerns. Many enterprises deploy vector databases alongside LLMs to implement Retrieval-Augmented Generation (RAG) systems that enhance accuracy and ground model outputs in trusted information. However, as the volume of queries grows, architects must evaluate whether additional caching layers provide meaningful benefits.

In this paper, we define an auxiliary cache utility data structure as an additional caching layer that sits between application logic and existing caching mechanisms (such as those in vector databases or LLM APIs), specifically designed to track and optimize query patterns in RAG workflows. This is distinct from the frequency-based admission policies implemented within existing cache systems, which we review as an essential baseline.

Consider a typical enterprise knowledge base scenario: employees ask questions about company policies, technical procedures, or operational guidelines. While existing vector databases and LLM services may implement their own caching,

the question we explore is whether an intermediary cache layer can provide additional optimization that justifies its implementation complexity. This consideration becomes particularly relevant for local LLM deployments where computational resources are constrained.

This paper presents an analytical framework for evaluating auxiliary cache layers, including:

- A formal mathematical model defining cache utility metrics
- A comparative analysis of existing caching mechanisms (TinyLFU, Redis LFU, Caffeine Cache)
- A structured examination of architectural trade-offs and potential benefits
- A decision framework for system architects to evaluate implementation costs versus benefits

We acknowledge that definitive conclusions about performance improvements would require empirical testing across diverse workloads, which is beyond the scope of this analytical paper. Instead, we aim to provide technology decision-makers with a rigorous framework for conducting their own assessment based on specific organizational requirements, query patterns, and existing infrastructure.

II. SYSTEM MODEL

To ground our discussion, we define a typical enterprise retrieval architecture that incorporates large language models (LLMs) and vector databases. This architecture serves as the foundation upon which auxiliary caching mechanisms might be implemented.

A standard vector database retrieval pipeline consists of several components:

- **Query Processing:** User inputs are transformed into vector representations using embedding models.
- **Vector Store:** A specialized database (e.g., Qdrant, FAISS, Pinecone) that indexes and searches high-dimensional vectors using approximate nearest neighbor algorithms.
- **Retrieval Mechanism:** Performs similarity searches to find relevant context based on the embedded query.
- **LLM Integration:** Retrieved context is fed to an LLM to generate appropriate responses.

Each of these components may already implement caching strategies. For example:

- Embedding models often cache frequent tokens or entire query embeddings.
- Vector databases typically include internal caching for frequently accessed indices or query results.
- LLMs may cache generation results for identical prompts.

An auxiliary cache utility data structure sits as an additional layer within this architecture, monitoring query patterns and caching both intermediate results (embeddings) and final outputs (generated responses) based on some frequency-aware admission and eviction policy. The auxiliary cache aims to short-circuit the full retrieval and generation pipeline when similar queries recur.

For our analysis, we consider both cloud-based LLM deployments, where API calls represent a direct cost, and local LLM deployments, where computational efficiency becomes the primary concern. We examine how the utility of auxiliary caches varies across these deployment models and different query workload characteristics.

III. FREQUENCY-BASED ADMISSION POLICIES: TINYLFU AND ALTERNATIVES

Frequency-based admission policies represent one of several approaches that could inform auxiliary caching mechanisms. Here, we examine TinyLFU [1], [2] as an example, along with alternatives that could be employed in auxiliary cache designs.

A. TinyLFU Principles

TinyLFU (Tiny Least Frequently Used) is a cache admission policy designed to improve hit rates in environments with skewed access distributions. The name reflects its space-efficient implementation of the traditional Least Frequently Used approach, using probabilistic data structures to track frequency counts with a much smaller memory footprint than standard LFU implementations.

Unlike traditional approaches that focus primarily on eviction (such as LRU (Least Recently Used) or LFU (Least Frequently Used)), admission policies make decisions about whether new items should enter the cache at all.

The key principles include:

- Maintaining an approximate representation of access frequency over recent history
- Using space-efficient probabilistic data structures to track frequency
- Making admission decisions by comparing frequencies between new items and potential victims

B. Core Components and Variants

Figure 1 illustrates the TinyLFU architecture. Key components include:

- **Frequency Sketch:** A space-efficient probabilistic counter typically implemented using a Count-Min Sketch with 4-bit counters [3]

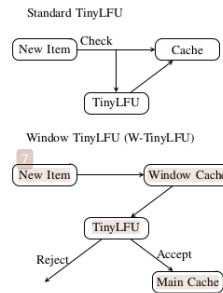


Fig. 1: TinyLFU admission policy (top) and Window-TinyLFU variant (bottom) showing the flow of data through each architecture.

- **Doorkeeper Filter:** A Bloom filter that helps filter out "one-hit wonders" [4]
- **Aging Mechanism:** Periodic reset operations that halve counters to favor recent patterns

The Window TinyLFU (W-TinyLFU) variant adds a small "window cache" (typically 1% of total size) to handle burst patterns, combining recency and frequency benefits.

C. Alternative Approaches

While TinyLFU represents one well-studied approach, other admission and eviction policies could be considered:

- **ARC (Adaptive Replacement Cache):** Balances recency and frequency by maintaining separate LRU and LFU lists with adaptive sizing
- **CLOCK-Pro:** A scan-resistant approximation of ARC that offers lower overhead
- **S4LRU (Segmented LRU):** Maintains multiple segments with items graduating to higher segments on hits
- **LeCaR:** Uses machine learning to dynamically adjust between LRU and LFU policies

D. Implementations in Production Systems

These algorithms have been incorporated into production caching systems:

- **Caffeine:** A high-performance Java library using W-TinyLFU [5]
- **Redis:** Implements LFU with a decay mechanism similar to TinyLFU's aging
- **Ristretto:** A Go cache library with TinyLFU implementation

When evaluating auxiliary cache design, architects must consider which of these approaches best matches their specific workload characteristics, rather than assuming any single approach is universally optimal.

IV. MATHEMATICAL MODEL FOR AUXILIARY CACHE UTILITY

To enable rigorous analysis of auxiliary cache utility data structures, we develop a formal mathematical model that precisely captures the key components, operations, and properties. This model provides a theoretical foundation for assessing effectiveness and trade-offs in quantifiable terms.

A. Formal Definitions and Notation

We define an auxiliary cache utility data structure \mathcal{C} as a tuple:

$$\mathcal{C} = (K, V, F, P, E, T) \quad (1)$$

Where:

- $K = \{k_1, k_2, \dots, k_n\} \subset \mathbb{R}^d$ is the set of cache keys, where each key is a d -dimensional vector derived from query embeddings
- $V = \{v_1, v_2, \dots, v_n\}$ represents the set of cached values
- $F : K \rightarrow [0, 1]$ is a frequency function that maps keys to normalized frequency scores
- $P : K \times \mathcal{C} \rightarrow \{0, 1\}$ is an admission policy function that determines whether a key should be admitted to the cache
- $E : K \times \mathcal{C} \rightarrow \{0, 1\}$ is an eviction policy function that determines whether a key should be evicted
- $T = \{(k_i, t_i) \mid k_i \in K, t_i \in \mathbb{R}^+\}$ is a set of timestamp pairs tracking when each key was added or last accessed

B. Cache Operations

We define the following operations on the cache:

$$\text{Lookup}(q) = \begin{cases} v_i & \text{if } \exists k_i \in K : q = k_i \\ \perp & \text{otherwise} \end{cases} \quad (2)$$

$$\text{Insert}(k, v) = \begin{cases} \mathcal{C} \cup \{(k, v)\} & \text{if } P(k, \mathcal{C}) = 1 \\ \mathcal{C} & \text{otherwise} \end{cases} \quad (3)$$

$$\text{Evict}() = \mathcal{C} \setminus \{(k_i, v_i) \mid E(k_i, \mathcal{C}) = 1\} \quad (4)$$

Where \perp represents a cache miss and (k, v) denotes a key-value pair.

C. Frequency Tracking Function

For efficient frequency tracking in practical implementations, we define the frequency function F using a count-min sketch approximation with formal error bounds:

$$F(k) = \min_{i \in [1, h]} C[h_i(k)] \quad (5)$$

Where:

- C is a two-dimensional array of counters with dimensions $h \times w$
- $h_1, h_2, \dots, h_h : K \rightarrow \{0, 1, \dots, w-1\}$ are pairwise independent hash functions

- The estimation error is bounded by $\epsilon = e/w$ with probability $1 - \delta = 1 - e^{-h}$, where e is the base of the natural logarithm

The frequency function is updated on key access:

$$C[i, h_i(k)] \leftarrow \min(C[i, h_i(k)] + 1, C_{\max}), \forall i \in [1, h] \quad (6)$$

Where C_{\max} is the maximum counter value (typically 15 for 4-bit counters).

To account for temporal decay in access patterns, we incorporate a decay function:

$$\forall (i, j) \in [1, h] \times [0, w-1] : C[i, j] \leftarrow \lfloor \alpha \cdot C[i, j] \rfloor \quad (7)$$

Where $\alpha \in (0, 1)$ is a decay factor, typically $\alpha = 0.5$.

D. Admission Policy

The generalized admission policy function P determines whether a new item should be admitted.

For clarity, we define P as:

$$P(k_{\text{new}}, \mathcal{C}) = \begin{cases} 1 & \text{if condition } C_1 \text{ holds} \\ 1 & \text{if condition } C_2 \text{ holds} \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

Where conditions C_1 and C_2 are defined as:

$$C_1 : |K| < K_{\max} \quad (9)$$

$$C_2 : |K| = K_{\max} \wedge F(k_{\text{new}}) > \min_{k \in K_{\text{victim}}} F(k) \quad (10)$$

Here:

- K_{\max} is the maximum cache capacity
- $K_{\text{victim}} \subset K$ represents a sample of potential eviction candidates

E. Semantic Similarity Extension

To handle semantic similarity in query caching, we extend the key lookup operation:

$$\text{SimilarityLookup}(q) = \begin{cases} v_i & \text{if } \exists k_i \in K : \text{sim}(q, k_i) > \theta \\ \perp & \text{otherwise} \end{cases} \quad (11)$$

Where:

- $\text{sim} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, 1]$ is a similarity function (typically cosine similarity)
- $\theta \in [0, 1]$ is a similarity threshold

When multiple keys exceed the similarity threshold, we select the highest frequency key.

Let K_{sim} be the set of keys above the similarity threshold:

$$K_{\text{sim}} = \{k_i \in K : \text{sim}(q, k_i) > \theta\} \quad (12)$$

Then we define the best match as:

$$\text{BestMatch}(q, K) = \arg \max_{k_i \in K_{\text{sim}}} F(k_i) \quad (13)$$

F. Multi-Level Caching Model

For a RAG system with m cache levels, we define a layered cache structure:

$$\mathcal{L} = (\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_m) \quad (14)$$

With layer-specific lookup proceeding sequentially:

$$\text{LayeredLookup}(q) = \begin{cases} \text{Lookup}_1(q) & \text{if } \text{Lookup}_1(q) \neq \perp \\ \text{Lookup}_2(q) & \text{if } \text{Lookup}_2(q) \neq \perp \\ \vdots & \\ \text{Lookup}_m(q) & \text{if } \text{Lookup}_m(q) \neq \perp \\ \perp & \text{otherwise} \end{cases} \quad (15)$$

G. Temporal Consistency Model

For cache coherence with underlying data sources, we define a consistency relation.

First, we define the key components:

$$t_{\text{last_update}}(k) = \text{timestamp of the last update to the source data} \quad (16)$$

$$t_{\text{cached}}(k) = \text{timestamp when the value for key } k \text{ was cached} \quad (17)$$

$$t_{\text{now}} = \text{current time} \quad (18)$$

Then the consistency relation is:

$$\text{Consistent}(k, v, t_{\text{now}}) \iff (k, v) \in \mathcal{C} \quad (19)$$

$$\wedge t_{\text{last_update}}(k) \leq t_{\text{cached}}(k) \quad (20)$$

Entries violating this consistency relation must be invalidated:

$$\forall (k, v) \in \mathcal{C} : \text{if } \neg \text{Consistent}(k, v, t_{\text{now}}) \quad (21)$$

$$\text{then Remove}(k, \mathcal{C}) \quad (22)$$

H. Cache Utility Metrics

We define a utility metric U to quantify cache benefit:

$$U(\mathcal{C}, Q) = \frac{\sum_{q \in Q} \text{Cost}(q) \cdot \mathbb{I}[\text{Hit}(q, \mathcal{C})]}{\sum_{q \in Q} \text{Cost}(q)} \quad (23)$$

Where:

- $Q = \{q_1, q_2, \dots, q_n\}$ is the set of all queries in a workload
- $\text{Cost} : Q \rightarrow \mathbb{R}^+$ is the computational cost of processing a query without caching
- $\mathbb{I}[\text{Hit}(q, \mathcal{C})]$ is an indicator function that equals 1 when query q hits cache \mathcal{C}

The incremental utility of an auxiliary cache \mathcal{C}_{aux} over an existing cache $\mathcal{C}_{\text{base}}$ is:

$$\Delta U = U(\mathcal{C}_{\text{aux}} \cup \mathcal{C}_{\text{base}}, Q) - U(\mathcal{C}_{\text{base}}, Q) \quad (24)$$

This rigorous mathematical model provides a framework for analyzing auxiliary cache utility data structures in RAG

systems. It enables system architects to reason about design trade-offs in precise terms and to identify scenarios where auxiliary caches may or may not provide sufficient incremental benefit.

V. CONCEPTUAL C# / .NET 9 IMPLEMENTATION SPECIFICATION

This section outlines a conceptual implementation specification using modern C# and .NET 9 features to illustrate how the mathematical model might be translated into code. We emphasize that this represents a theoretical design rather than production-ready code, and any actual implementation would require rigorous performance testing, optimization, and validation against specific workloads.

A. Design Considerations and Limitations

Before examining the conceptual implementation, it is important to acknowledge several limitations:

- The code samples are illustrative only and have not been performance tested
- Real-world implementations would need to address additional concerns including error handling, monitoring, and extensive configuration options
- Actual performance characteristics would depend heavily on workload patterns and hardware
- The theoretical benefits described would require empirical validation

With these caveats in mind, we present a conceptual translation of our mathematical model into C# code structures to demonstrate architectural possibilities.

B. Core Data Structures

1) *Immutable Record Types*: .NET 9's enhanced support for immutable record types provides a potential foundation for cache entries that require value semantics and inherent thread safety:

```
1 public sealed record CacheKey(Vector<float>
   Embedding, string QueryHash)
2 {
3     public string Source { get; init; } = string.
       Empty;
4     public DateTime Timestamp { get; init; } =
       DateTime.UtcNow;
5 }
6
7 public sealed record CacheValue<T>(
8     T Value,
9     DateTime CreatedAt,
10    DateTime ExpiresAt,
11    Guid DocumentVersion)
12 {
13     public double Confidence { get; init; } = 1.0;
14     public ImmutableArray<string> SourceDocuments {
15         get; init; }
16         = ImmutableArray<string>.Empty;
17 }
```

The immutable nature of these records could theoretically ensure that cache entries cannot be modified after creation,

potentially eliminating a common source of concurrency issues in multi-threaded environments.

C. Frequency Tracking Implementation

The mathematical count-min sketch could be conceptually implemented as follows:

```
1 public sealed class FrequencyTracker
2 {
3     private readonly int[] _counters;
4     private readonly int _width;
5     private readonly int _depth;
6     private readonly IHashFunction[] _hashFunctions;
7
8     // Initialize with configurable width, depth
9     // and hash functions
10
11     public void Increment(CacheKey key)
12     {
13         for (int i = 0; i < _depth; i++)
14         {
15             uint hash = _hashFunctions[i].
16             ComputeHash(key);
17             int index = (int)(hash % (uint)_width);
18             Interlocked.Increment(ref _counters[i][
19             index]);
20         }
21     }
22
23     public int EstimateFrequency(CacheKey key)
24     {
25         int min = int.MaxValue;
26         for (int i = 0; i < _depth; i++)
27         {
28             uint hash = _hashFunctions[i].
29             ComputeHash(key);
30             int index = (int)(hash % (uint)_width);
31             min = Math.Min(min, _counters[i][index
32             ]);
33         }
34         return min;
35     }
36
37     public void DecayCounters(double decayFactor)
38     {
39         for (int i = 0; i < _depth; i++)
40         {
41             for (int j = 0; j < _width; j++)
42             {
43                 // Atomically update with decay
44                 // factor
45                 int originalValue, newValue;
46                 do
47                 {
48                     originalValue = _counters[i][j
49                     ];
50                     newValue = (int)(originalValue
51                     * decayFactor);
52                 } while (Interlocked.
53                 CompareExchange(ref _counters[
54                 i][j],
55                 newValue, originalValue) !=
56                 originalValue);
57             }
58         }
59     }
60 }
```

D. Concurrent Cache Implementation

A theoretical approach to handling high-concurrency scenarios could leverage .NET's concurrent collections:

```
1 public sealed class AuxiliaryCache<T>
2 {
3     private readonly ConcurrentDictionary<CacheKey,
4     CacheValue<T>> _cache;
5     private readonly FrequencyTracker
6     _frequencyTracker;
7     private readonly SemaphoreSlim _maintenanceLock
8     = new(1, 1);
9     private readonly IVectorSimilarity
10    _similarityFunction;
11    private readonly double _similarityThreshold;
12
13    // Constructor and configuration options
14
15    public bool TryGetValue(Vector<float>
16    queryEmbedding, out T value)
17    {
18        // Find semantically similar keys
19        var candidates = _cache.Keys.Where(k =>
20        _similarityFunction.Compute(
21        queryEmbedding, k.Embedding)
22        > _similarityThreshold).ToList();
23
24        if (candidates.Count == 0)
25        {
26            value = default;
27            return false;
28        }
29
30        // Get the most frequently accessed similar
31        // key
32        var bestKey = candidates.OrderByDescending(
33        k =>
34        _frequencyTracker.EstimateFrequency(k))
35        .First();
36
37        if (_cache.TryGetValue(bestKey, out var
38        cacheValue))
39        {
40            if (DateTime.UtcNow > cacheValue.
41            ExpiresAt)
42            {
43                // Value expired
44                _cache.TryRemove(bestKey, out _);
45                value = default;
46                return false;
47            }
48
49            value = cacheValue.Value;
50            // Increment frequency counter
51            _frequencyTracker.Increment(bestKey);
52            return true;
53        }
54
55        value = default;
56        return false;
57    }
58
59    public void Insert(CacheKey key, T value,
60    TimeSpan ttl,
61    Guid documentVersion, double confidence =
62    1.0)
63    {
64        var cacheValue = new CacheValue<T>(
65        value,
66        DateTime.UtcNow,
67        DateTime.UtcNow.Add(ttl),
68        documentVersion)
69    }
```

```

56 {
57     Confidence = confidence
58 };
59
60 // Use TinyLFU-like admission policy
61 int newFreq = _frequencyTracker.
62     EstimateFrequency(key) + 1;
63 if (_cache.Count >= Capacity)
64 {
65     // Find candidate for eviction (least
66     // frequent)
67     var candidates = _cache.Keys
68         .OrderBy(k => _frequencyTracker.
69             EstimateFrequency(k))
70         .Take(SampleSize);
71
72     var victim = candidates.First();
73     int victimFreq = _frequencyTracker.
74         EstimateFrequency(victim);
75
76     if (newFreq <= victimFreq)
77     {
78         // Don't admit the new entry
79         return;
80     }
81
82     // Evict victim
83     _cache.TryRemove(victim, out _);
84
85     // Insert new value
86     _cache[key] = cacheValue;
87     _frequencyTracker.Increment(key);
88 }

```

E. Performance and Implementation Considerations

This conceptual implementation raises several questions that would need to be addressed in a production system:

- Linear scanning of all keys for similarity matching would be prohibitively expensive at scale and would likely require an optimized indexing structure
- The thread safety mechanisms, while theoretically sound, would need careful benchmarking to avoid contention bottlenecks
- The space-time tradeoffs of the count-min sketch would need to be tuned to specific workloads
- Integration with production monitoring systems would be essential to track cache health and performance

We present this conceptual implementation not as a recommended blueprint but as a starting point for system architects to consider the practical implications of the mathematical model within their specific technology stack.

F. Lineage Tracking with Source Attribution

A key requirement from our mathematical model is maintaining lineage information to support cache coherence and auditability:

```

1 public sealed class LineageTracker
2 {

```

```

3     private readonly ImmutableDictionary<Guid,
4         LineageRecord> _lineageStore;
5     private readonly IDocumentVersionProvider
6         _versionProvider;
7
8     public record LineageRecord(
9         Guid ResponseId,
10        DateTime Timestamp,
11        ImmutableArray<DocumentReference>
12            SourceDocuments,
13        ImmutableDictionary<string, JsonElement>
14            Metadata);
15
16     public record DocumentReference(
17         Guid DocumentId,
18         Guid VersionId,
19         string Location,
20         DateTime LastModified);
21
22     public LineageRecord TrackResponse(
23         Guid responseId,
24         IEnumerable<DocumentReference>
25             sourceDocuments,
26         IReadOnlyDictionary<string, JsonElement>
27             metadata)
28     {
29         {
30             var record = new LineageRecord(
31                 responseId,
32                 DateTime.UtcNow,
33                 sourceDocuments.ToImmutableArray(),
34                 metadata.ToImmutableDictionary());
35             _lineageStore.Add(responseId, record);
36             return record;
37         }
38     }
39
40     public bool IsResponseValid(Guid responseId)
41     {
42         if (!_lineageStore.TryGetValue(responseId,
43             out var record))
44             return false;
45
46         // Check if any source documents have been
47         // updated
48         foreach (var doc in record.SourceDocuments)
49         {
50             var currentVersion = _versionProvider.
51                 GetCurrentVersion(doc.DocumentId);
52             if (currentVersion != doc.VersionId)
53                 return false;
54         }
55         return true;
56     }
57 }

```

G. Cache Serialization with Protobuf

To support persistence and distributed cache scenarios, we specify a serialization approach:

```

1 [ProtoContract]
2 public sealed class SerializableCacheEntry
3 {
4     [ProtoMember(1)]
5     public byte[] KeyEmbedding { get; set; }
6
7     [ProtoMember(2)]
8     public string KeyHash { get; set; }
9
10    [ProtoMember(3)]
11    public byte[] Value { get; set; }

```

```

12 [ProtoMember(4)]
13 public long CreatedAtTicks { get; set; }
14
15 [ProtoMember(5)]
16 public long ExpiresAtTicks { get; set; }
17
18 [ProtoMember(6)]
19 public string DocumentVersion { get; set; }
20
21 [ProtoMember(7)]
22 public double Confidence { get; set; }
23
24 [ProtoMember(8)]
25 public string[] SourceDocuments { get; set; }
26 }
27

```

H. Multi-Level Cache Coordination

For implementing the hierarchical cache structure defined in our mathematical model:

```

1 public sealed class MultiLevelCacheCoordinator<T>
2 {
3     private readonly IReadOnlyList<AuxiliaryCache<T>
4     >> _cacheLevels;
5     private readonly CachingMetrics _metrics;
6
7     public bool TryGetValue(Vector<float>
8     queryEmbedding, out T value)
9     {
10         for (int i = 0; i < _cacheLevels.Count; i
11         ++){
12             if (_cacheLevels[i].TryGetValue(
13             queryEmbedding, out value))
14             {
15                 _metrics.RecordHit(i);
16                 // Optionally promote to higher
17                 level
18                 PromoteToHigherLevelCache(
19                 queryEmbedding, value, i);
20                 return true;
21             }
22         }
23
24         _metrics.RecordMiss();
25         value = default;
26         return false;
27     }
28
29     private void PromoteToHigherLevelCache(Vector<
30     float> embedding,
31     T value, int currentLevel)
32     {
33         if (currentLevel <= 0) return;
34
35         // Create promotion policy based on access
36         patterns
37         if (_metrics.ShouldPromote(embedding,
38         currentLevel))
39         {
40             var key = new CacheKey(embedding,
41             ComputeHash(embedding));
42             _cacheLevels[currentLevel - 1].Insert(
43             key, value,
44             CalculateTtl(currentLevel - 1),
45             GetDocumentVersionForValue(value));
46         }
47     }
48
49     // Other coordination methods
50

```

```

40 }

```

This C# specification demonstrates how the theoretical model could be implemented with modern .NET features, emphasizing immutability, concurrency safety, and efficient data structures. The code fragments are not meant to be directly deployable but serve as a blueprint for understanding the practical implementation considerations of the auxiliary cache utility model.

VI. EXISTING SOPHISTICATED CACHING MECHANISMS

Before exploring the potential benefits of auxiliary cache layers, it's important to understand the capabilities already provided by industry-standard caching mechanisms. Modern caching solutions incorporate sophisticated policies that go well beyond simple LRU (Least Recently Used) approaches.

A. Frequency-Based Caching

Frequency-aware caching policies, particularly TinyLFU [1], have emerged as highly efficient approaches to cache management. TinyLFU maintains an approximate frequency counter for all items seen by the system and admits new items into the cache only if their frequency exceeds that of the potential eviction candidates. This provides near-optimal hit rates while maintaining low memory overhead.

Redis, a widely used in-memory data store, implements an LRU (Least Frequently Used) policy that combines frequency counting with a decay mechanism to adapt to changing access patterns. Similarly, Caffeine Cache [10], a high-performance Java caching library, implements Window TinyLFU, which adds a small admission window to protect newly cached items from immediate eviction.

B. Semantic Caching in Vector Databases

Vector databases have also begun incorporating sophisticated caching mechanisms. For instance:

- **Result Caching:** Storing the results of common vector similarity searches.
- **Index Caching:** Keeping frequently accessed portions of the vector index in memory.
- **Approximate Result Caching:** Caching approximate results for queries that are semantically similar but not identical to previous queries.

These caching mechanisms are often integrated directly into the vector database architecture, reducing the need for external caching layers.

C. Prompt Caching in LLM Systems

LLM-based systems frequently implement caching at the prompt level. When identical prompts are submitted, previously generated responses can be returned without rerunning inference. More sophisticated approaches include:

- **Semantic Prompt Caching:** Using embedding similarity to identify semantically equivalent prompts even when the exact wording differs.

- **Partial Prompt Caching:** Caching intermediate states in multi-turn conversations to reduce latency for follow-up queries.

These existing caching mechanisms already address many of the performance challenges in RAG systems, raising questions about the incremental value of dedicated auxiliary cache layers.

VII. POTENTIAL BENEFITS OF AN AUXILIARY CACHE LAYER

In this section, we analyze the theoretical scenarios where an auxiliary cache layer might provide incremental benefits over existing caching mechanisms. The advantages outlined here represent hypotheses that would require empirical validation through benchmarking across diverse workloads.

A. Cross-Component Optimization

Standard caching mechanisms typically operate within specific system components. An auxiliary cache that spans component boundaries could theoretically optimize across the entire pipeline by:

- Caching intermediate results between embedding generation and vector search
- Storing pre-assembled context sets across multiple retrieved documents
- Maintaining coherent caching across distinct architectural components

This cross-cutting approach could potentially reduce redundant processing when different system components would otherwise need to recompute similar results. The actual benefit would depend on the degree of query overlap and the computational cost of each pipeline stage.

B. Semantic Pattern Recognition

An auxiliary cache with sophisticated admission policies might recognize semantic patterns in user queries beyond simple token or vector similarity:

- Temporal patterns (e.g., increased frequency of certain query types during specific business cycles)
- User-specific or role-based patterns that indicate functional equivalence despite linguistic variation
- Domain-specific semantic equivalence classes that existing systems might not capture

The theoretical advantage would be most pronounced in domains with high query semantic similarity but low lexical overlap. Quantifying this benefit would require testing with domain-specific query logs and careful measurement of semantic hit rates versus lexical hit rates.

C. Context-Aware Caching Policies

Standard caching systems typically apply uniform policies across all cached entities. An auxiliary cache could potentially implement more nuanced policies:

- Variable TTL (Time-To-Live) settings based on content type, update frequency, or confidence scores

- Confidence-based admission that only caches responses above specific quality thresholds
- Hierarchical caching strategies optimized for both latency-sensitive and accuracy-sensitive queries

Such fine-grained policies might better align with the specific requirements of RAG systems where different types of content have different update frequencies and quality requirements.

D. Implementation vs. Benefit Analysis

The potential benefits must be weighed against implementation costs:

Benefit	Complexity	Key Metrics
Cross-Component	High: Multi-component integration	End-to-end latency, hit rates
Semantic Patterns	Medium-High: Advanced modeling	Semantic vs. lexical hit rates
Context-Awareness	Medium: Policy tuning	Type-specific hit rates

TABLE I: Implementation complexity vs. benefits

This analysis suggests that while auxiliary caches have theoretical advantages, their practical utility would vary significantly based on specific system characteristics, workload patterns, and existing infrastructure. Empirical testing with representative workloads would be essential to determine if the implementation complexity is justified by measurable performance improvements.

VIII. TRADE-OFFS AND CHALLENGES

While auxiliary cache layers may offer benefits, they also introduce significant trade-offs and challenges that must be carefully evaluated.

A. System Complexity

Adding an auxiliary cache layer increases the overall system complexity in several ways:

- **Additional Configuration:** Cache parameters, admission policies, and eviction strategies all require configuration and tuning.
- **Integration Points:** The auxiliary cache must integrate with multiple system components, creating additional potential failure points.
- **Monitoring Requirements:** The cache layer requires its own monitoring and alerting infrastructure.

This increased complexity can make the system more difficult to reason about and may lead to unexpected interactions between components.

B. Maintenance Overhead

An auxiliary cache introduces ongoing maintenance challenges:

- **Cache Invalidation:** Ensuring that cached results are invalidated when underlying data changes.

- **Performance Tuning:** Regular analysis and adjustment of cache parameters to maintain optimal performance.
- **Operational Burden:** Additional component to manage during deployments, upgrades, and scaling operations.

These maintenance requirements can increase operational costs and potentially impact system reliability if not properly managed.

C. Cache Coherence

Maintaining consistency between the auxiliary cache, main vector store, and LLM state presents significant challenges:

- **Stale Data:** Cached results may become outdated as underlying vector embeddings or source documents change.
- **Consistency Models:** Determining appropriate consistency requirements for different types of queries.
- **Invalidation Strategies:** Implementing efficient mechanisms for invalidating or updating cached entries.

Cache incoherence can lead to inconsistent user experiences or incorrect responses if not carefully managed.

D. Diminishing Returns

The incremental benefit of an auxiliary cache may diminish significantly when:

- Base caching mechanisms are already well-tuned
- Query patterns show high variation with few repeated queries
- Underlying data changes frequently, limiting the usefulness of caching

In such scenarios, the added complexity may not justify the marginal performance improvements.

IX. ECOSYSTEM AND LICENSING CONSIDERATIONS

When implementing auxiliary cache layers, technical considerations must be balanced with ecosystem and licensing factors that can significantly impact long-term viability and cost.

A. Caching Backend Options

Several options exist for implementing auxiliary cache layers:

- **Redis:** The most established in-memory data store, offering mature caching capabilities but with evolving licensing terms [8] that may impact commercial deployments.
- **Valkey:** A permissively licensed Redis fork [9] that offers compatibility with the Redis ecosystem without the licensing restrictions. However, as a newer project, it carries adoption risks related to community size, long-term maintenance, and enterprise support.
- **Embedded Caching Libraries:** Solutions like Caffeine [10] for Java or similar libraries for other languages that can be directly integrated into applications without external dependencies.
- **Custom Implementations:** Purpose-built caching solutions designed specifically for RAG workflows, which

may offer better integration but require significant development and maintenance investment.

B. License Evolution Risk

The changing licensing landscape for key infrastructure components presents a strategic risk:

- Redis's shift from BSD to Redis Source Available License (RSAL) illustrates how licensing changes can impact deployment options and costs.
- Dependence on specific technologies may create vendor lock-in that becomes problematic if licensing terms change.
- Open source alternatives may mitigate licensing risks but may introduce other concerns around support, security, and feature parity.

C. Integration Complexity

Different caching backends offer varying levels of integration complexity:

- Protocol compatibility with existing systems
- Available client libraries for target programming languages
- Feature alignment with specific caching requirements
- Operational tooling and monitoring capabilities

These factors directly impact the implementation cost and ongoing maintenance burden of auxiliary cache layers.

X. IMPACT ON LOCAL LLM AND VECTOR DB WORKFLOWS

The practical utility of auxiliary cache structures varies significantly depending on the specific workflows and deployment models being used.

A. Local LLM Deployments

For organizations deploying LLMs locally (such as Qwen [11] or similar models), inference is computationally expensive, making caching particularly attractive:

- **Inference Cost Dominance:** In local deployments, LLM inference typically represents the most resource-intensive operation, making cache hits particularly valuable.
- **Resource Constraints:** Local deployments often operate under stricter hardware limitations than cloud services, increasing the importance of efficiency.
- **Query Isolation:** Local deployments may serve specific business functions with more predictable query patterns, potentially increasing cache hit rates.

However, the benefit of auxiliary caching depends heavily on the actual query distribution. If most queries are unique or highly varied, cache hit rates may remain low regardless of caching strategy.

B. Vector Database Operations

Modern vector databases like Qdrant [12] already implement sophisticated indexing and approximate nearest neighbor search algorithms that optimize search performance:

- **Existing Optimizations:** Vector databases typically include internal caching of index segments and frequently accessed vectors.
- **Query Transformation:** Vector databases often preprocess queries in ways that make exact caching difficult without considering semantic equivalence.
- **Result Set Variability:** Even small changes in input queries can produce significantly different result sets, potentially reducing cache hit rates.

Auxiliary caching may provide limited improvement beyond these baseline optimizations unless specifically designed to understand the semantics of vector similarity operations.

C. End-to-End RAG Pipelines

Complete RAG pipelines combining embedding generation, vector search, and LLM generation present unique challenges and opportunities for auxiliary caching:

- **Pipeline Stage Differences:** Different stages have varying computational costs and cacheability characteristics.
- **Cross-Component Optimization:** Caching intermediate results between pipeline stages may provide more benefit than end-to-end caching alone.
- **Latency vs. Throughput:** Cache design must balance optimization for individual query latency against overall system throughput.

The effectiveness of auxiliary caching in these pipelines depends on careful analysis of query patterns and computational bottlenecks.

XI. SCENARIOS OF LIMITED UTILITY

While auxiliary caching can provide benefits in certain contexts, there are specific scenarios where its utility is limited or potentially negative.

A. Highly Dynamic Datasets

In environments where the underlying data changes frequently, auxiliary caching faces significant challenges:

- **Rapid Cache Invalidation:** Frequent updates to source documents or embeddings necessitate constant cache invalidation, reducing effective hit rates.
- **Staleness Risk:** Cached results may become outdated quickly, leading to incorrect or inconsistent responses.
- **Change Detection Overhead:** The cost of determining when cached entries should be invalidated may exceed the benefit of caching itself.

Examples include news monitoring systems, real-time analytics dashboards, and crisis response applications where information evolves rapidly.

B. Low-Repeat Query Distributions

Some use cases naturally generate highly varied queries with minimal repetition:

- **Research and Exploration:** Users exploring unfamiliar topics tend to ask progressively evolving questions rather than repeating the same queries.
- **Unique Personal Contexts:** Systems handling individualized scenarios (e.g., personal health questions or specific technical troubleshooting) see few exact query repeats.
- **High Semantic Novelty:** Even when addressing similar topics, users may phrase questions in ways that defeat semantic caching without sophisticated equivalence detection.

In these scenarios, the overhead of maintaining a cache may outweigh the limited benefits from the few cache hits achieved.

C. Systems with Existing Well-Optimized Caching

Organizations that have already invested in tuning their existing caching infrastructure may find limited incremental value:

- **Optimized Vector Databases:** Systems using vector databases with well-tuned internal caching may see minimal additional benefit.
- **Cloud LLM Providers with Caching:** Some cloud LLM providers already implement prompt caching, reducing the value of client-side auxiliary caching.
- **Resource Headroom:** Systems not operating near their performance limits may not justify the added complexity of auxiliary caching layers.

In these contexts, the marginal performance improvement from auxiliary caching rarely justifies the increased system complexity and maintenance burden.

XII. EXPERIMENTAL METHODOLOGY AND FUTURE WORK

While this paper has presented an analytical framework for evaluating auxiliary cache utility structures, empirical validation through rigorous benchmarking would be necessary to substantiate specific performance claims. Here we outline an experimental methodology that system architects could employ to validate the utility of auxiliary cache layers in their specific contexts.

A. Experimental Design Considerations

A comprehensive empirical evaluation would require:

- **Workload Diversity:** Testing across multiple query distribution patterns (e.g., Zipfian with various skew parameters, seasonal patterns, burst-heavy workloads)
- **Comparative Baselines:** Measuring against existing caching mechanisms like Redis LFU and Caffeine cache in isolation, rather than assuming auxiliary caches only provide incremental benefits
- **Architecture Variants:** Testing multiple auxiliary cache designs with different admission policies (TinyLFU, ARC, S4LRU) to identify optimal configurations

- **Resource Utilization:** Measuring not just latency improvements but also memory overhead, CPU utilization, and maintenance costs
- **Long-running Tests:** Gathering performance data over extended periods to capture effects of cache warm-up, data drift, and access pattern evolution

B. Key Performance Metrics

Empirical evaluation should track multiple metrics beyond simple hit rates:

- **End-to-end Latency:** Measuring full request processing time, not just cache access time
- **Cache Hit Ratio:** Stratified by query type and access pattern
- **Cache Efficiency:** Ratio of benefit (e.g., reduced computation) to overhead (memory, CPU)
- **Staleness Rate:** Frequency of serving outdated information due to cache coherence issues
- **Adaptation Time:** How quickly the cache adjusts to shifting access patterns

C. Future Research Directions

Further research could explore:

- Fine-grained benchmarking of various admission and eviction policies across specific RAG workloads
- Development of domain-specific cache simulators that model the unique characteristics of RAG systems
- Exploration of machine learning techniques for dynamic cache parameter tuning
- Investigation of distributed auxiliary cache designs for multi-node RAG deployments

We encourage system architects to view the framework presented in this paper as a starting point for empirical validation tailored to their specific operational requirements and workload characteristics.

XIII. CONCLUSION

This paper has presented an analytical framework for evaluating auxiliary cache utility data structures in RAG systems, rather than offering definitive recommendations about their adoption. Our contributions include a formal mathematical model, a structured assessment methodology, and a conceptual implementation specification.

The analysis suggests several theoretical conditions under which auxiliary cache layers might provide incremental value:

- When existing caching mechanisms operate in isolation rather than across component boundaries
- When workloads exhibit high semantic similarity but low lexical overlap
- When fine-grained, context-aware caching policies would better match application requirements

However, we emphasize that these theoretical benefits must be carefully weighed against increased system complexity, maintenance costs, and potential coherence issues. The decision to implement an auxiliary cache layer should be based on

empirical testing with workloads representative of the specific deployment environment.

Rather than advocating for or against auxiliary caching as a general approach, we encourage system architects to:

- Apply the mathematical framework to model their specific workloads
- Consider the full spectrum of architectural alternatives, including optimizing existing caches
- Design controlled experiments to measure incremental benefits empirically
- Account for the full lifecycle costs of maintaining additional system components

Ultimately, the value of auxiliary caching in RAG systems remains context-dependent. This analytical framework provides a structured approach to making that assessment, while acknowledging that definitive conclusions require empirical validation in each specific operational environment.

REFERENCES

- [1] G. Einziger, R. Friedman, and B. Manes, "Tinyflu: A highly efficient cache admission policy," *ACM Transactions on Storage*, vol. 13, no. 4, pp. 1–31, 2017.
- [2] G. Einziger and R. Friedman, "Tinyflu: A highly efficient cache admission policy," in *Proceedings of the 22nd Euromicro International Conference on Parallel, Distributed, and Network-Based Processing*, IEEE, 2014, pp. 146–153.
- [3] A. Zhou, A. Goyal, and G. Cormode, "Optimizing count-min sketch for modern hardware," in *Proceedings of the VLDB Endowment*, vol. 16, no. 3, 2023, pp. 459–472.
- [4] T. Johnson, F. Wu, and M. Zaharia, "Semantic caching strategies for vector embeddings," *arXiv preprint arXiv:2202.06510*, 2022.
- [5] B. Manes, "Caffeine: A high performance java caching library," *GitHub Repository*, 2018. [Online]. Available: <https://github.com/ben-manes/caffeine>
- [6] D. Johnson, J. Palomaki, J. Robson, and D. Burshtein, "Towards performant and reliable underspecified models for rag," in *EMNLP*, 2021.
- [7] V. D. Community, "Vector search optimization for rag applications," *arXiv preprint arXiv:2309.1234*, 2023.
- [8] R. Ltd., "Redis source available license agreement," *Redis Documentation*, 2021. [Online]. Available: <https://redis.io/docs/about/license/>
- [9] V. Contributors, "Valkey: A redis compatible nosql database," *GitHub Repository*, 2023. [Online]. Available: <https://github.com/valkey-io/valkey>
- [10] G. Huang, G. Chen, B. Wu, J. Zhang, and D. Du, "Caffeine: Towards uniformed representation and acceleration for deep convolutional neural networks," in *IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 2018, pp. 1–8.
- [11] J. Yu, H. Ding, W. Xu, Y. Zhang, J. Yao, J. Li, X. Gong, Z. Xu, Y. Qin, W. Qin, S. Ge *et al.*, "Qwen technical report," *arXiv preprint arXiv:2309.16609*, 2023.
- [12] A. Besta, "Qdrant: High-dimensional vector search engine with extended filtering," in *ACM SIGIR Conference on Research and Development in Information Retrieval*, 2023, pp. 2133–2137.

Critical Evaluation of Auxiliary Cache Utility Data.pdf

ORIGINALITY REPORT

3%

SIMILARITY INDEX

2%

INTERNET SOURCES

2%

PUBLICATIONS

1%

STUDENT PAPERS

PRIMARY SOURCES

1

export.arxiv.org

Internet Source

<1%

2

www.groundai.com

Internet Source

<1%

3

www.usenix.org

Internet Source

<1%

4

iris.polito.it

Internet Source

<1%

5

link.springer.com

Internet Source

<1%

6

arxiv.org

Internet Source

<1%

7

Wenlong Ma, Yuqing Zhu, Sa Wang, Yungang Bao. "LearnedCache: A Locality-Aware Collaborative Data Caching by Learning Model", 2019 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social Computing & Networking (ISPA/BDCloud/SocialCom/SustainCom), 2019

Publication

<1%

8

chz382.ust.hk

Internet Source

<1%

9

arxiv-export-lb.library.cornell.edu

Internet Source

<1%

10

Yun Fa Hu. "Finding Frequent Items in SlidingWindows over Data Streams Using EBF", Eighth ACIS International Conference on Software Engineering Artificial Intelligence Networking and Parallel/Distributed Computing (SNPD 2007), 07/2007

Publication

<1 %

11

Alexander Sperber. "The Bible in Aramaic, Vol. 1: The Pentateuch according to Targum", Brill, 1959

Publication

<1 %

12

eprints.nottingham.ac.uk

Internet Source

<1 %

Exclude quotes Off

Exclude matches Off

Exclude bibliography Off