

# MeanCache: User-Centric Semantic Caching for LLM Web Services

Waris Gill<sup>1</sup>, Mohamed Elidrisi<sup>2</sup>, Pallavi Kalapatapu<sup>2</sup>, Ammar Ahmed<sup>4</sup>, Ali Anwar<sup>4</sup>, Muhammad Ali Gulzar<sup>1</sup>

<sup>1</sup>Virginia Tech <sup>2</sup>Cisco <sup>3</sup>University of Minnesota

waris@vt.edu, {melidris, pkalapat}@cisco.com, {ahme0599, aanwar}@umn.edu, gulzar@cs.vt.edu

**Abstract**—Large Language Models (LLMs) like ChatGPT and Llama have revolutionized natural language processing and search engine dynamics. However, these models incur exceptionally high computational costs. For instance, GPT-3 consists of 175 billion parameters, where inference demands billions of floating-point operations. Caching is a natural solution to reduce LLM inference costs on repeated queries. However, existing caching methods are incapable of finding semantic similarities among LLM queries nor do they operate effectively on contextual queries, leading to unacceptable *false hit-and-miss* rates.

This paper introduces MeanCache, a user-centric semantic cache for LLM-based services that identifies semantically similar queries to determine cache hit or miss. Using MeanCache, the response to a user’s semantically similar query can be retrieved from a local cache rather than re-querying the LLM, thus reducing costs, service provider load, and environmental impact. MeanCache leverages Federated Learning (FL) to collaboratively train a query similarity model without violating user privacy. By placing a local cache in each user’s device and using FL, MeanCache reduces the latency, costs, and enhances model performance, resulting in lower false-hit rates. MeanCache also encodes context chains for every cached query, offering a simple yet highly effective mechanism to discern contextual query responses from standalone queries. Our experiments benchmarked against the state-of-the-art caching method reveal that MeanCache attains an approximately 17% higher F-score and a 20% increase in precision during semantic cache hit-and-miss decisions while performing even better on contextual queries. It also reduces the storage requirement by 83% and accelerates semantic cache hit-and-miss decisions by 11%.

**Index Terms**—Large Language Models, Semantic Cache, Embedding, Contextual Queries, Cache, Privacy-Preserving AI

## I. INTRODUCTION

Large Language Models (LLMs) like ChatGPT [1], Google Bard [2], Claude [3], and Llama [4] have demonstrated remarkable capabilities in understanding and generating human language, leading to significant advancements in applications ranging from search engines to conversational agents. LLMs are increasingly integrated into platforms like the Perplexity AI search engine, Rabbit OS [5], and Arc browser [6].

**Motivation.** Generating responses to user queries with LLMs, such as GPT-3, requires substantial computations and poses environmental challenges [7], [8], [9], [10]. For example, GPT-3’s 175B parameters in float16 format consume 326 GB memory, exceeding single GPU capacities and necessitating multi-GPU deployments [7]. These requirements lead to high operational costs. Consequently, LLM-based services charge users and limit query rates [11], [12]. Prior studies have observed that users frequently submit similar queries to web

services [13], [14], [15] (approximately 33% of search engine queries being resubmitted [15]), suggesting opportunities for optimization by avoiding redundant computations.

Caching serves as an effective technique in traditional web services to address duplicate search queries, avoiding redundant computations, significantly improving response time, reducing the load on query processors, and enhancing network bandwidth utilization [15], [13], [16], [17]. If applicable to LLMs-based web services, such caching can substantially impact billions of floating point operations, thereby decreasing operational costs and environmental impact.

**Problem.** Existing caching techniques [18], [15], [14], [13] use keyword matching, which often struggles to capture the semantic similarity among similar queries to LLM-based web services, resulting in a significantly low hit rate. For instance, existing caches do not detect the semantic similarity between “How can I increase the battery life of my smartphone?” and “Tips for extending the duration of my phone’s power source”, leading to a cache miss. Recently, Zhu et al. [19] and GPTCache [20] present server-side semantic caching for the LLMs-based services to address the limitations of keyword-matching caching techniques. If a new query is semantically similar to any query in a cache, the server returns the response from the cache. Otherwise, a model multiplexer selects the most suitable LLM for the query to generate the response.

Existing semantic caches have several limitations. First, they demand a significantly large central cache to store the queries and responses of all users, which is unscalable and violates users’ privacy. Second, they incur the network cost of sending a user query to the server even if there is a cache hit. An end user will still be charged for the query even if the query is served from the server cache. Third, they use a single server-side embedding model to find the semantic similarity among queries, which does not generalize to each user’s querying patterns. For instance, Google Keyboard [21] adapts to each user’s unique writing style and embeds such personalized behaviors to enhance the accuracy of the next word prediction model. Fourth, they employ Llama-2 to enhance the accuracy of semantic matching [20]; however, in practice, such models perform billions of operations to generate embeddings, offsetting the benefits of the cache. Lastly, they are only effective on standalone queries, resulting in unbearably high false hit rates for contextually different but semantically similar queries.

**Key Insights and Contributions of MeanCache.** This work introduces a novel *user-centric* semantic caching system

called MeanCache. MeanCache provides a privacy-preserving caching system that returns the response to similar queries directly from the user’s local cache, bypassing the need to query the LLM-based web service. MeanCache achieves these goals in the following ways.

To address privacy concerns associated with central server-side caching, MeanCache introduces a user-side cache design ensuring that the user’s queries and responses are never stored outside of the user’s device. To find a semantic match between a new query and cached queries, MeanCache uses smaller embedding models such as MPNet [22] to generate embeddings for semantic matching locally. Previous work has shown that a smaller model can achieve performance comparable to larger models on custom tasks [23], [24], [25].

Due to different contexts around queries, LLM may return different responses for semantically similar queries. For such queries, MeanCache also records the contextual chain, parent queries already in the cache, for a given query. To find a response for a contextual query, MeanCache verifies the context of a contextual query by matching a given query’s context with the cached query’s context chain to accurately retrieve responses to contextual queries.

Each user may not have sufficient queries to customize an embedding model that can help find a semantic match between new queries and cached queries. To address this, MeanCache utilizes federated learning (FL), which exploits data silos on user devices for private training for collaborative learning [26], [27], [28], [29], [?], thereby personalizing an embedding model for each user. This privacy-preserving training not only customizes the embedding model to the user’s querying patterns but also enhances the performance (*i.e.*, accuracy) of semantic caching without compromising user privacy (*i.e.*, without storing user data on the web server).

The runtime performance of MeanCache is primarily influenced by the time taken to match a new query embedding vector with existing ones in the cache to find a semantically similar query. The search time is directly proportional to the dimensions of the embedding vector. To optimize runtime performance, MeanCache compresses the embedding vector by leveraging principal component analysis (PCA) [30], [31], [32], effectively reducing the size of the embedding vector (*i.e.*, projecting it to lower dimensional space). MeanCache also offers an adaptive cosine similarity threshold, which is also collaboratively computed using FL, to improve accuracy in finding semantic matches between queries.

**Evaluations.** We compare MeanCache with GPTCache [20]. GPTCache is closely related to MeanCache and has received over 6,000 stars on GitHub [33]. We benchmark MeanCache’s performance against GPTCache on the GPTCache’s dataset [34] to demonstrate its effectiveness and highlight the improvements MeanCache offers over existing solutions.

MeanCache surpasses GPTCache [20] by achieving a 17% higher F-score and approximately a 20% increase in precision in end-to-end deployment for identifying duplicate queries to LLM-based web services. MeanCache’s performance on contextual queries is even more impressive when compared

to GPTCache (baseline). For contextual queries, MeanCache achieves a 25% higher F-score and accuracy, and a 32% higher precision over the baseline. MeanCache’s embedding compression utility approximately reduces storage and memory needs by 83% and results in 11% faster semantic matching while still outperforming the state-of-the-art GPTCache.

**Artifact Availability:** MeanCache is implemented in the Flower FL framework [35]. The complete source code and contextual queries dataset will be available at <https://github.com/SEED-VT/MeanCache>.

## II. BACKGROUND

**Federated Learning (FL).** FL is a distributed, privacy-preserving ML model training approach [36], [27], [37], [38]. In each FL round, a central server distributes a global model to participating clients, who train it on their local data and return the updated models. The server then aggregates these models to create a new global model for the next round. While several aggregation algorithms exist (FedAvg [28], FedProx [39], and FedMA [40]), FedAvg is most common, using the equation:

$$W_{global}^{t+1} = \sum_{k=1}^K \frac{n_k}{n} w_{k,t} \quad (1)$$

where  $W_{global}^{t+1}$  represents the new global model,  $w_{k,t}$  is the model of the  $k^{th}$  client at round  $t$ ,  $n_k$  is the sample count at the  $k^{th}$  client, and  $n$  is the total sample count across participating clients. This process continues for multiple rounds until convergence.

**Transformer and Embeddings.** The transformer architecture is based on the attention mechanism [41]. Transformers are extensively utilized in natural language processing (NLP) tasks such as machine translation, text summarization, and question-answering. The attention mechanism allows transformers to effectively capture long-range dependencies within sequences, while positional encoding explicitly incorporates information about the order of tokens. Transformers convert text into embeddings, representing words or sentences as dense vectors in high-dimensional space. These embeddings capture semantic meanings, ensuring that semantically similar words or sentences have similar vector representations [42], [43], [44]. Cosine similarity is a widely used metric for measuring the similarity between embeddings, with values ranging from -1 (completely opposite) to 1 (identical). For instance, embeddings for words such as *cat* and *dog* typically exhibit higher cosine similarity compared to *cat* and *car*, as the former pair is more closely related semantically. For embedding vectors  $\mathbf{E}_1$  and  $\mathbf{E}_2$ , cosine similarity is calculated as:

$$\text{cosine\_similarity} = \frac{\mathbf{E}_1 \cdot \mathbf{E}_2}{\|\mathbf{E}_1\| \cdot \|\mathbf{E}_2\|} \quad (2)$$

where  $\mathbf{E}_1 \cdot \mathbf{E}_2$  represents the dot product and  $\|\mathbf{E}_1\|$  and  $\|\mathbf{E}_2\|$  denote the vector magnitudes.

**Contextual Queries.** Interaction between end-users and LLM services typically involves two primary types of queries: standalone queries (e.g., Q1 *Draw a line in Python?*) and follow-

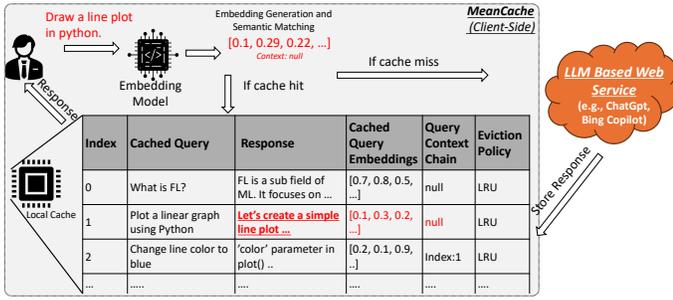


Fig. 1: MeanCache's Workflow.

up, or contextual queries (e.g., Q2 *Change the color to red*). Standalone queries can be resolved independently without any additional information, whereas contextual queries require prior context to provide accurate and relevant responses.

Consider a scenario in which both queries Q1 and Q2, along with their responses, have been cached. If a user later submits a new standalone query, Q3 *Draw a circle?*, followed by a contextual query, Q4 *Change the color to red*, Q4 might appear semantically similar to the cached query Q2. However, despite their similarity, Q4 references a different context (Q3 rather than Q1), thus requiring a distinct interpretation and response. Consequently, a semantic cache without effective context detection may yield incorrect cache hits.

### III. MEANCACHE'S DESIGN

MeanCache is a user-centric semantic cache optimized for user-side operation. Figure 1 illustrates the workflow of MeanCache for similar queries. Algorithm 1 further explains MeanCache's querying and population process programmatically. When a user submits a query to an LLM web service with MeanCache enabled, MeanCache computes the query embeddings (Line 1). These embeddings are then matched with the embedding of the cached queries using cosine similarity (Line 2). For every similar query found within the cache, MeanCache analyzes the context chain for every query and matches it with the conversational history of the submitted query (Line 4). If MeanCache finds a similar query with a similar context chain, the response is retrieved from the local cache and returned to the user (Line 7). Otherwise, MeanCache forwards the query to the LLM service to obtain the response. The query, its response, and embeddings are then stored in the cache (Line 9).

MeanCache harnesses the collective intelligence of multiple users to train a semantic similarity model, and its user-centric design addresses privacy and scalability issues. To achieve these, MeanCache takes the following design decisions. It employs a small embedding model with lower computational overhead than LLM based embedding models. It uses FL for collaborative training to fine-tune the embedding model without ever storing user data on a central server. This approach generates high-quality embeddings and improves the accuracy of embedding matching for retrieving similar queries. To handle contextual queries, MeanCache includes contextual chain information in its cache against every query

### Algorithm 1: MeanCache

---

**Input:** User query  $Q$ , User Query Context  $C_q$  (e.g., *null* if no parent)

**Output:** Response  $R$

- 1  $E_Q \leftarrow \text{encode}(Q)$  // compute the embedding of the query
- 2  $\text{similar\_queries} \leftarrow \text{FindSimilarQueriesInCache}(E_Q)$   
// retrieve top- $k$  similar cached queries
- 3  $\text{context\_match} \leftarrow \text{False}$  // flag to indicate if suitable context is found
- 4 **foreach** context  $C_i \in \text{similar\_queries}$  **do**
- 5     **if**  $C_i$  matches with  $C_q$  **then**
- 6          $\text{context\_match} \leftarrow \text{True}$
- 7         Retrieve response  $R$  from cache for  $C_i$
- 8 **if not**  $\text{context\_match}$  **then**
- 9      $R \leftarrow \text{LLMResponseAndEnrollInCache}(Q, E_Q, C_q)$   
// generate response and cache it
- 10 **return**  $R$

---

to identify if the cached response for a query is only applicable under a specific context. This design is capable of handling contextual chains. To reduce storage and memory overhead and expedite search time for finding similar queries in the cache, MeanCache compresses the embeddings using PCA.

#### A. FL Based Embedding Model Training

GPTCache [20] suggests using Llama to generate superior embeddings, thereby enhancing semantic matching accuracy. However, this approach has several limitations. LLMs not only sizable, being gigabytes in size, but they also require substantial computational resources to generate embeddings. Deploying such models, especially at the end-user level, is impractical due to their size and significant computational overhead for semantic matching. MeanCache employs a compact embedding model, which has a lower computational overhead compared to large embedding models. The smaller model may not provide the same level of accuracy as an LLM. However, a smaller model trained on customized tasks can match the performance of an LLM [23], [24], [25]. One challenge in using smaller embedding models is that each user may not have sufficient data to train and customize the embedding model.

MeanCache utilizes FL to exploit the vast amount of distributed data available on users' devices to train and personalize the smaller embedding model. FL allows the users to train the embedding model locally and learn the optimal threshold for cosine similarity. The updated weights and local threshold are shared with the server. The server aggregates the updated weights and cosine similarity threshold from multiple users to update the global model, which is redistributed back to the users. This approach ensures that the user's privacy is maintained, and the collective intelligence of multiple users is leveraged to improve the performance of the caching system. Figure 2 shows the overview of privacy-preserving training of the embedding model with FL.

In the first step, the server sends the initial weights of the embedding model ( $W_{global}^{t+1}$ ) and global threshold ( $\tau$ ) to a subset of users as shown in step 1 in Figure 2. The subset

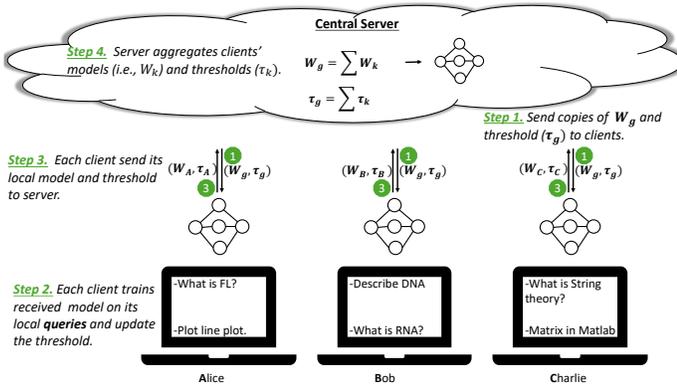


Fig. 2: Privacy-preserving model training in MeanCache.

of users is usually randomly selected or selected based on their battery level, network bandwidth, or performance history. Additionally, the server also sends the hyperparameters (e.g., learning rate, batch size, epochs) necessary for FL training of the embedding model.

1) *Client Training*: Upon receiving the embedding model ( $W_{global}^{t+1}$ ) from the server, each client replaces its local embedding model weights with the newly received weights. Subsequently, each client trains the embedding model locally on its unique local dataset (step 2 in Figure 2). To generate high-quality embeddings from unique and similar queries, MeanCache's clients training employs a multitask learning approach, integrating two distinct loss functions: contrastive loss [42] and multiple-negatives ranking loss [45], [42]. These loss functions update the weights of the embedding model during the training process. The contrastive loss function operates by distancing unique (non-duplicate) queries from each other within the embedding space, thereby facilitating the differentiation between duplicate and non-duplicate queries. Unlike contrastive loss, the multiple-negatives ranking loss function minimizes the distance between positive pairs (duplicate queries) amidst a large set of potential candidates i.e., multiple-negatives ranking loss does not concentrate on distancing unique queries and its objective is to draw positive pairs (similar queries) closer within the embedding space.

This learning approach enables MeanCache to adjust to diverse query patterns exhibited by users. For instance, some users may generate more repetitive queries compared to others, while certain users may not produce any repetitive queries at all. Interestingly, MeanCache's multitask learning objective can benefit from learning even from a user with no repetitive queries. This is because MeanCache's global embedding model ( $W_{global}^{t+1}$ ) will learn to widen the distance between unique queries, thereby effectively learning the true misses of the non-duplicate queries and minimizing the false hits during the search process. True miss happens when a similar query is not present in the cache. A false hit is when a query is found and returned from the cache, which is not actually similar.

2) *Finding the Optimal Threshold for Cosine Similarity*: After generating query embeddings using an embedding model, a similarity metric such as cosine similarity is used

to determine if the new query embeddings match the cached embeddings of past queries. This process involves setting a threshold for cosine similarity, which is a delicate balance.

In addition to privacy-preserving training of the embedding model, MeanCache also learns the optimal threshold ( $\tau$ ) for cosine similarity. The range of  $\tau$  is between 0 and 1. This threshold ( $\tau$ ) dictates the level of similarity above which a cached query is considered relevant to the current user query. Setting the threshold too low could result in numerous false hits, leading to retrieving irrelevant queries from the cache. Conversely, a threshold set too high might cause many false misses, where relevant queries are not retrieved from the cache.

During the client's local training, MeanCache determines this optimal threshold ( $\tau$ ) from the client's feedback to the cache query response. Even after receiving a cached response, a user requests a response from the LLM, MeanCache considers it as a false positive and adjusts its threshold. MeanCache varies the threshold  $\tau$  to find the optimal threshold that optimizes the F-score of the cache (Section IV-F). By finding the optimal threshold, MeanCache effectively balances between true hits and true misses, therefore yielding improved accuracy in semantic similarity matching to return the response from the cache on duplicate queries.

3) *Aggregation*: After client local training and finding the optimal threshold, each client sends updated weights of the global model ( $W_{global}^{t+1}$ ) and optimal threshold ( $\tau$ ) to the server (step 3 in Figure 2). The server aggregates the updated weights from multiple users to form a new embedding model ( $W_{global}^{t+1}$ ) using FedAvg [28] as shown in step 4 of Figure 2. The server also computes the mean of the received optimal thresholds from the clients for a global optimal threshold ( $\tau_{global}$ ). The benefit of finding  $\tau_{global}$  is that when a new user joins the system, the user will not have queries to find its own optimal threshold. In such cases, the system can use  $\tau_{global}$  as a starting point for semantic similarity.

After finding the global optimal threshold and the global embedding model, the server then redistributes the updated embedding model to the users for the next round of FL training. This process is repeated over several FL rounds to improve the semantic matching accuracy (i.e., lower false hits and false misses). After the completion of the FL training, each client will have access to a fine-tuned embedding model ( $W_{global}^{t+1}$ ) to generate high-quality embeddings that can capture the complex semantics of a user query.

4) *Embeddings Compression using PCA*: The substantial size of the embedding vector (e.g., Llama embeddings with a dimension of 4096) can lead to considerable overhead during the matching process of new query embeddings with cached queries embeddings. This is due to the search time being directly proportional to the dimensions of the embedding vector. Furthermore, high-dimensional embeddings demand more memory and storage. For example, the embeddings generated by Llama for a single query require approximately 32.05 KB of memory storage. Consequently, calculating the cosine similarity between two high-dimensional embeddings, specifically new query embeddings and each embedding in the cache,

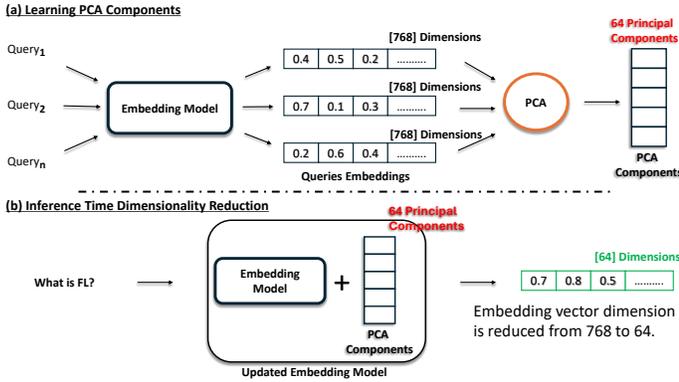


Fig. 3: Embeddings Compression using PCA.

becomes computationally demanding and time-intensive. To improve the search time for identifying similar queries in the cache and to reduce the client’s storage needs, it is essential to diminish the dimensionality of the embeddings while ensuring minimal impact on the MeanCache’s performance.

PCA is a dimensionality reduction technique that is widely used to compress high-dimensional data into a lower-dimensional space [30], [31], [32], while still maintaining the most important information. First, MeanCache generates embeddings for all the users’ queries using the embedding model. Next, MeanCache applies PCA to learn the principal components of all the queries embeddings generated in the previous step, as shown in Figure 3-a. MeanCache integrates the learned principal components as an additional layer in the embedding model. This new layer will project the original embeddings onto the lower dimensional space, producing compressed embeddings (Figure 3-b).

When a non-duplicate query is received, MeanCache uses the updated embedding model (with PCA layer) to generate the compressed embeddings (Figure 3-b) for the new query and store the query, response, and the compressed embedding in the cache. Storing the compressed embeddings in the cache will significantly reduce the storage and memory overhead of the embeddings. Next, when a duplicate query is received, MeanCache uses the same embedding model with PCA components to generate the compressed embeddings for this duplicate query and find similar queries in the cache. Since the embeddings are compressed, the search time for finding similar queries in the cache will be significantly reduced.

### B. Cache Population and MeanCache Implementation

Once the embedding model is trained within MeanCache on each user, it is deployed as depicted in Figure 1. Initially, when a new user starts using MeanCache, the local cache is vacant. During these interactions, if a user query’s response is not found in the cache, the request is forwarded to the LLM web service to retrieve the response, which is then inserted in the cache. If MeanCache finds semantically similar queries in the cache for any of the following queries from the user, it analyzes the context chain for every similar cached query and matches its embedding with the conversational history

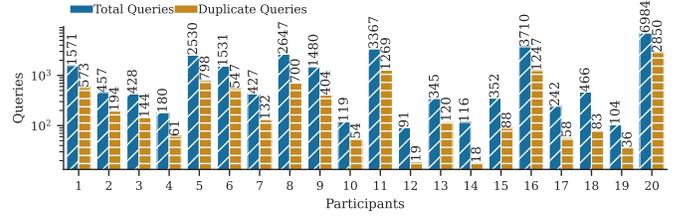


Fig. 4: Analysis of real-world ChatGPT conversations reveals that, on average, 31% of queries are similar to previously submitted queries.

of the submitted query. If MeanCache finds a similar query with a similar context chain, the response is retrieved from the local cache and returned to the user. Otherwise, MeanCache forwards the query to the LLM web service to obtain the response. The query and its response and embeddings are then stored in the cache.

MeanCache is a python-based application that is built on the Flower FL framework [35]. A user can submit LLM queries via this application to take advantage of the local cache. The central server, which orchestrates the FL training, may reside on the LLM web service. We employ the Sbert [42] library to train MPNet [22] and Albert [46] on each client and to generate query embeddings. To efficiently execute a cosine similarity search between a query embedding and cached embeddings, we utilize Sbert’s semantic search, which can handle up to 1 million entries in the cache. MeanCache cache storage is persistent and built using DiskCache [47] library.

### C. Validating the Prevalence of Similar Queries

We conducted a privacy-preserving study with 20 ChatGPT users, analyzing over 27K queries. Our participants included university professors, developers, and graduate students who are regular ChatGPT users and experienced in running Python scripts. We provide these participants detailed instructions to download their queries and responses, run analysis scripts locally, and share aggregated results (e.g., total and duplicate queries, field/profession). Individual queries are not shared, and thus, the conversations remain private. We find that about 31% of user queries were similar to previous ones, suggesting a user-centric cache can reduce LLM inference costs (Figure 4). As this study was conducted in an academic setting, the ratio of repeated queries may vary in other contexts.

## IV. EVALUATION

Our evaluation answers the following research questions.

- How does MeanCache perform in comparison to baseline in terms of performance metrics (§IV-B)?
- How accurately does MeanCache retrieve contextual queries from the cache (§IV-C)?
- Is it possible to reduce the embedding dimension to save storage space and accelerate semantic search while outperforming the baseline (§IV-D)?
- Is it possible to train an embedding model in a privacy-preserving manner without the centralized data (§IV-E)?

| Metrics   | Standalone Queries |                    |                     | Contextual Queries |             |
|-----------|--------------------|--------------------|---------------------|--------------------|-------------|
|           | GPT-Cache          | Mean-Cache (MPNet) | Mean-Cache (Albert) | GPT-Cache          | MeanCache   |
| F score   | 0.56               | 0.73               | 0.68                | 0.67               | <b>0.93</b> |
| Precision | 0.52               | 0.72               | 0.66                | 0.66               | <b>0.98</b> |
| Recall    | 0.85               | 0.78               | 0.77                | 0.71               | 0.79        |
| Accuracy  | 0.72               | 0.85               | 0.81                | 0.61               | <b>0.86</b> |

TABLE I: MeanCache outperforms GPTCache (baseline) on both standalone and contextual queries.

- What effect does the cosine similarity threshold have on the performance of MeanCache (§IV-F)?
- GPTCache [20] suggests using Llama 2 to generate embeddings to improve semantic matching. Is it feasible to use Llama 2 to compute embeddings at the user side for semantic matching (§IV-G)?

#### A. Evaluation Settings

We conduct evaluations of MeanCache against the baseline [20] to demonstrate that MeanCache achieves optimal performance while preserving user-privacy (*i.e.*, without storing the user queries at the server). For a fair comparison between MeanCache and baseline, we employ the optimal configuration as described in the GPTCache study [20]. This configuration utilizes Albert [46] and sets the cosine similarity threshold at 0.7 to determine the cache hit or miss.

1) *Transformer Models and Datasets*: For extensive evaluations of MeanCache, we utilize the Llama 2 [4], MPNet [22], and Albert [46] transformer models to generate embeddings.

We evaluate MeanCache using the GPTCache dataset. The dataset is partitioned into training, testing, and validation subsets. The training and validation datasets are randomly distributed among the clients, each receiving non-overlapping data points. During local training, each client utilizes its training dataset to update its local embedding model and employs the validation dataset to determine the optimal threshold for cosine similarity (Section IV-F). The testing dataset, located at the central server, facilitates a fair comparison between MeanCache and GPTCache [20]. Since there does not exist any dataset of contextual queries, we generate a synthetic dataset using GPT-4 consisting of 450 queries, including duplicates, non-duplicates, and contextual queries, to evaluate MeanCache performance on contextual queries.

2) *Experimental Setup*: The experiments are conducted on a high-performance computing cluster, equipped with 128 cores, 504 GB of memory, and four A100 Nvidia GPUs, each with 80 GB of memory. We utilize the Flower FL [35] library to simulate a FL setup. Additionally, the SBERT [42] library is employed to train the embedding model on each client. The number of clients participating in FL training are 20. The number of clients is restricted due to the limited size of the GPTCache dataset, which is inadequate for distribution among hundreds of clients. However, we believe MeanCache results are not influenced by the number of clients, and the evaluation setup of MeanCache is consistent with the evaluation standard in FL [48], [49].

3) *Evaluation Metrics*: In caching systems, the efficacy has traditionally been gauged by cache hit-and-miss rates. A cache hit implies the data or query is retrieved from the cache, whereas a cache miss indicates the opposite. Semantic caching introduces a nuanced classification: true and false hits, alongside true and false misses. A *true hit* signifies a correct match between a query and a similar cached query, whereas a *false hit* is an incorrect match with a non-similar cached query. A *true miss* signifies when a query does not have a similar cached query, whereas a *false miss* happens when a query has a similar cached query but is not returned from the cache. Thus, traditional hit/miss metrics are potentially misleading in semantic caches. For example, a query might incorrectly match with an irrelevant cached query (deemed a hit traditionally) due to semantic matching. We adopt precision, recall, F score, and accuracy for a comprehensive evaluation of MeanCache against baseline. These metrics are defined as follows:

**Precision.** The ratio of true positive hits to all positive hits (including both true positives and false positives). In semantic caching, this measures how many of the queries matched to a cached query are correctly matched. Precision =  $\frac{TP}{TP+FP}$  where  $TP$  represents true positives (true hits) and  $FP$  represents false positives (false hits).

**Recall.** The ratio of true positive hits to all relevant items (including both true positives and false negatives). In semantic caching, this assesses the proportion of correctly matched queries out of all queries that should have been matched to a cached query. Recall =  $\frac{TP}{TP+FN}$  where  $FN$  represents false negatives (false misses).

**$F_\beta$  Score.** A weighted harmonic mean of precision and recall, balancing the two based on the value of  $\beta$ .  $\beta > 1$  gives more weight to recall, while  $\beta < 1$  emphasizes precision.

$$F_\beta = (1 + \beta^2) \cdot \frac{\text{Precision} \times \text{Recall}}{(\beta^2 \times \text{Precision}) + \text{Recall}}$$

**Accuracy.** The ratio of correctly identified queries (both true hits and true misses) to all queries. Accuracy =  $\frac{TP+TN}{TP+TN+FP+FN}$  where  $TN$  represents true negatives (true misses).

#### B. MeanCache Comparison with Baseline

We evaluate MeanCache against GPTCache to assess improvements in precision, recall, and F score. MeanCache FL model training is discussed in Section IV-E, and the optimal threshold is covered in Section IV-F.

We select a sample of 1000 queries, 30% of which are repeated queries (*i.e.*, 300 queries are repeated), and load these 1000 queries as cached queries. Note that repeated queries are usually fewer than non-repeated queries. Thus, we use 30% as repeated queries, a similar percentage previously observed for web services [15].

Initially, we send a new set of one thousand queries to Llama 2 (*i.e.*, without any semantic cache) to establish a baseline for response times. We limit responses to 50 tokens to reflect practical response sizes, although actual sizes can be much larger. Note that MeanCache’s performance is not

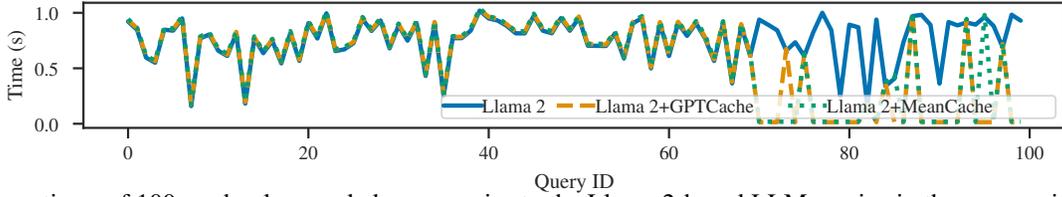


Fig. 5: Response times of 100 randomly sampled user queries to the Llama 2-based LLM service in three scenarios: without any semantic cache, with GPTCache, and with MeanCache. The integration of a semantic cache does not add significant overhead to non-duplicate queries, meaning it does not impede the performance of the LLM-based service. Moreover, it significantly reduces the average response times for duplicate queries (70-99) by serving them from the local cache.

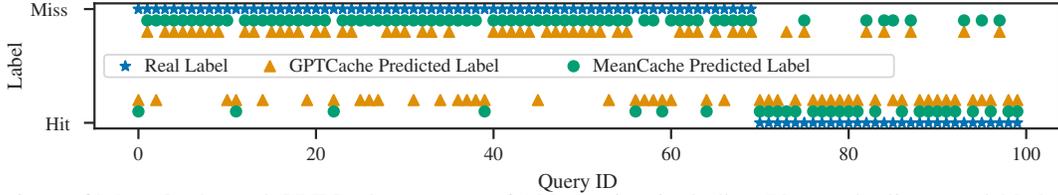


Fig. 6: Comparison of MeanCache and GPTCache on a set of 100 queries, including 70 non-duplicate and 30 duplicate queries, sent to the Llama 2-based LLM service. Queries 0 to 69 are non-duplicate (*i.e.*, *real label is miss*), and GPTCache produces significantly higher false hits on these unique queries compared to MeanCache.

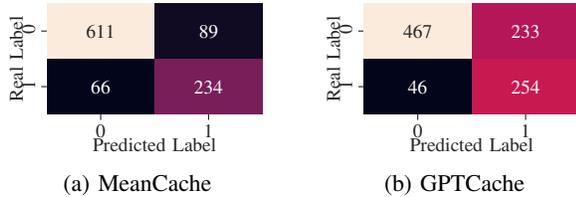
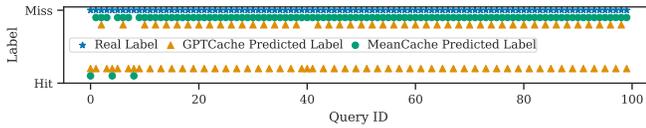


Fig. 7: Confusion matrices for MeanCache and GPTCache on 1000 queries comprising 700 unique and 300 duplicate queries. Among the 700 unique queries, MeanCache produces only 89 false hits, while GPTCache generates 233 false hits.

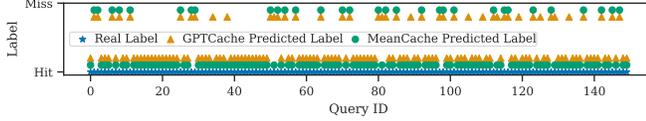
dependent on the response as it only matches the queries. Next, we send these queries to Llama 2 based local LLM service with MeanCache and GPTCache to measure the response times and performance metrics (*e.g.*, precision, recall, F-score), respectively. An analysis of a random subset of 100 queries (70 non-duplicate and 30 duplicate) from the 1000 queries shows the impact of caching on response times in Figure 5 and cache hit and miss rates in Figure 6. While MeanCache is evaluated on 1,000 standalone and 450 contextual queries (Table I), we use this smaller subset of 100 queries solely for visualization purposes, as displaying all queries would create cluttered figures. The y-axis in Figure 5 shows the response time of each query without any cache (Llama 2), with GPTCache (Llama 2 + GPTCache), and with MeanCache (Llama 2 + MeanCache). In Figure 6, x-axis represents the query and the y-axis represents each semantic cache hit and miss alongside the real label. Figure 5 demonstrates that implementing a semantic cache does not impede the performance (queries ranging from 0 to 69) and improves the user experience as response times reduce on duplicate queries (queries 70 to 99).

However, Figure 6 shows that GPTCache produces significant false hits on non-duplicate queries (queries 0 to 69) compared to MeanCache. Each false hit means the user receives an incorrect cached response, requiring them to resend the query to the LLM service, leading to a poor user experience. To prioritize precision, we adjust the beta value in the F score to 0.5, valuing precision twice as highly as recall to ensure user satisfaction by avoiding false positives. This decision is driven by the need to minimize user inconvenience caused by incorrect cache hits. A false miss (low recall) does not require user intervention as the false miss query will be automatically routed to the LLM. Thus, precision in semantic caching is more important than recall.

Table I and the confusion matrices in Figure 7 highlight MeanCache’s superior performance over GPTCache on the 1000 user queries. Notably, MeanCache with MPNet achieves a precision of 0.72, significantly surpassing GPTCache’s 0.52. This superiority is evident in the lower false positive rates (*i.e.*, false hits) shown in Figure 7a. The number of false hits for MeanCache is 89 (Figure 7a), while GPTCache has 233 false hits, as depicted in Figure 7b. In practical terms, this means that with GPTCache, the end user has to manually resend 233 queries to the LLM service to get the correct responses, compared to only 89 queries with MeanCache. While GPTCache’s recall is higher than MeanCache’s, as we discussed earlier, precision is significantly more important than recall, and MeanCache outperforms GPTCache in this regard. Overall, the F-score of MeanCache with MPNet is 0.73 and 0.68 with the Albert embedding model, both of which outperform GPTCache’s F-score of 0.56.



(a) Ideally, all 100 queries should result in misses. However, GPTCache incorrectly produces 54 *false hits*, while MeanCache yields only 3.



(b) MeanCache yields 8% more *true hits* than the baseline.

Fig. 8: Performance on Contextual Queries: MeanCache vs. Baseline. (a) reports MeanCache’s fewer false hits 3 vs. 54 of GPTCache. (b) reports higher true hits by MeanCache.

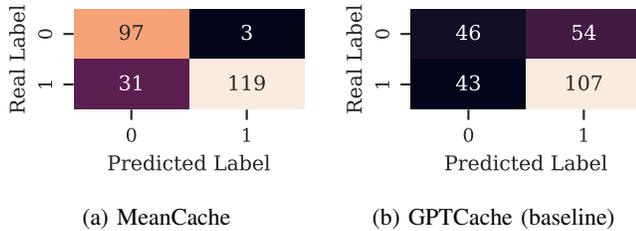


Fig. 9: For contextual queries, MeanCache reports three false hits, compared to 54 false hits by GPTCache.

**Summary.** In an end-to-end deployment, MeanCache significantly outperforms GPTCache. It demonstrates a 17% higher F score and a 20% increase in precision in optimal configuration. The substantial reduction in false cache hits enhances the end-user experience.

### C. Contextual Queries

Section II describes contextual queries. GPTCache, lacking the capability to detect contextual queries, incorrectly identifies such queries as cache hits, resulting in inaccurate responses. MeanCache addresses this limitation by verifying the context chain (Algorithm 1) of semantically matched queries. This verification ensures that contextual queries (*e.g.*, Q4 in Section II) correctly yield a *true cache miss*, thus prompting an appropriate request to the LLM service.

On the dataset of 450 contextual queries (see Section IV-A1), first, we populate the cache (MeanCache and baseline) with 200 queries (100 standalone and 100 contextual queries of the standalone queries). Next, we send 150 duplicate queries (75 standalone + 75 contextual) and 100 non-duplicate queries (a total of 250 queries) to the cache-enabled LLM. Figure 8 shows the true label (whether the query should be returned from the cache or not) and the corresponding GPTCache and MeanCache performance (predicted label). Note that in Figure 8a, all the queries should be answered by the LLM; in other words, there should be no hits. However, GPTCache has 54 false hits, while MeanCache has only three

false hits. This is also shown in the confusion matrix in Figure 9. Table I (Column-3) summarizes the comparative results. MeanCache outperforms GPTCache by over 25% in both F-Score and accuracy. Additionally, MeanCache achieves 32% higher precision compared to the baseline.

**Summary.** GPTCache’s low performance stems from its inability to consider contextual information, leading to high false hit rates. MeanCache demonstrates superior handling of contextual queries, with a 25% improvement in accuracy over the baseline.

### D. Embedding Compression and Impact on Storage Space

Clients often face storage limitations compared to web servers. Storing embeddings in the local cache on the user side for semantic search demands memory storage. Various models yield embeddings with differing vector sizes; for example, the MPNet and Albert models produce an output embedding vector of 768 dimensions, whereas the Llama 2 model’s embeddings dimension size is 4096. The embedding vector size also affects semantic search speeds, where smaller vectors could enhance speed and lower resource demands.

Figure 10 illustrates the effects of MeanCache dimension reduction utility on the storage, semantic matching speed (overhead), and MeanCache’s performance (F score). The x-axis indicates the number of queries stored in the cache, while the y-axis shows storage size, average search time, and the F score in respective graphs. MeanCache-Compressed (MPNet) and MeanCache-Compressed (Albert) represent instances where MeanCache decreases the embedding dimensions from 768 to 64 by employing the compression, as detailed in MeanCache design (Section III).

Figure 10 demonstrates that increased stored queries linearly raise storage needs. Yet MeanCache with compression enabled drastically lowers storage needs by 83% compared to GPTCache. Figure 10 also indicates that compression decreases the average search time, with MeanCache enabled compression approximately 11% faster. Moreover, despite a slight decrease in F score with compression enabled, MeanCache still surpasses GPTCache. Furthermore, given the evidence from Section IV-E (Figures 11 and 12) that MPNet produces more precise embeddings, and it is also clear from Figure 10 that MPNet’s embeddings are particularly resilient to compression and excel in semantic matching.

**Summary.** The application of embedding compression optimization in MeanCache offers substantial benefits, including an 83% savings in storage and an 11% faster search process, while still outperforming the established baseline (GPTCache).

### E. Privacy Preserving Embeddings Model Training

Storing clients queries on the server side presents a potential privacy risk. To address this, each client can retain its local data on its own device. The ensuing challenge is

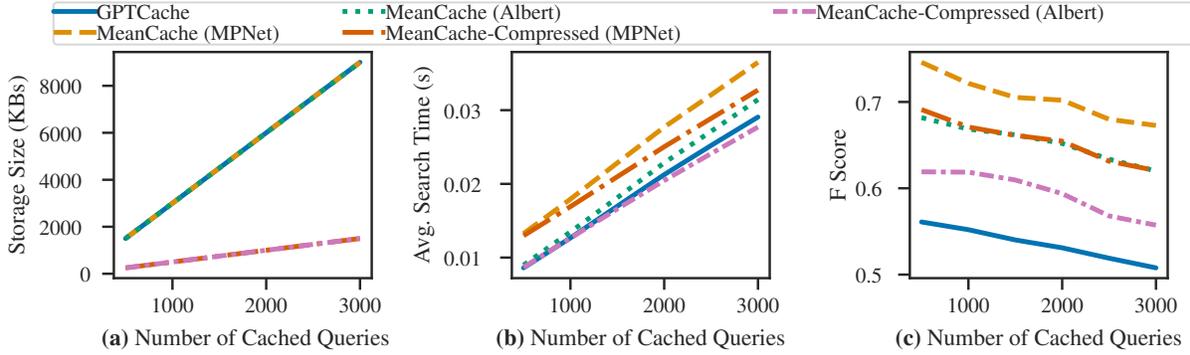


Fig. 10: (a) MeanCache’s embedding compression reduces storage by 83% compared to GPTCache. (b) MeanCache’s semantic matching speed is 11% faster with compression enabled, while still outperforming GPTCache. (c) MeanCache’s F score is slightly lower with compression enabled, but it still outperforms GPTCache.

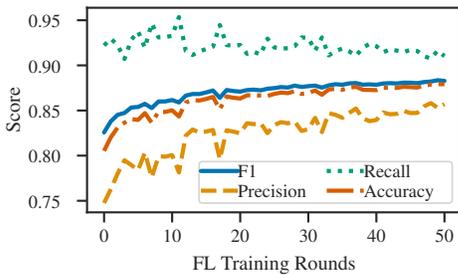


Fig. 11: *MPNet*’s FL training helps generate high-quality embeddings.

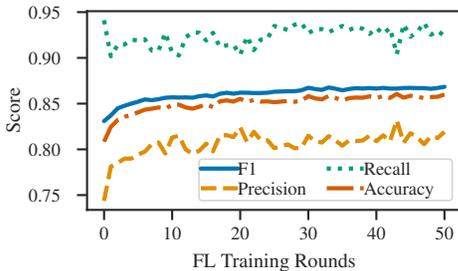


Fig. 12: FL training boosts MeanCache’s query matching precision with *Albert*.

how to train an embedding model that can also utilize the distributed data from all clients. FL is recognized for training neural networks in a privacy-preserving manner. As such, MeanCache employs FL to train and fine-tune an embedding model, thereby preserving privacy and leveraging the dataset residing on the client’s side. In this section, our objective is to evaluate whether FL training can progressively enhance the embedding model to generate high-quality embeddings for user queries. To simulate this scenario, we distribute the training dataset among 20 clients. In each round, we sample 4 clients, conducting a total of 50 FL training rounds. Each client trains its embedding model for 6 epochs, operating on a

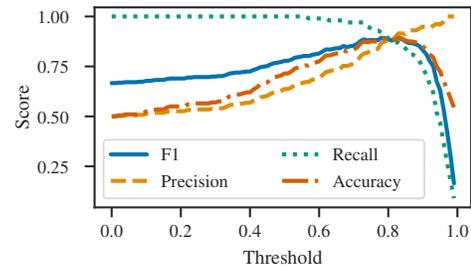


Fig. 13: MeanCache optimizes *MPNet*’s threshold.

dedicated A100 GPU. We conduct two experiments with the Albert and MPNet models. The batch size is set to 256 for the Albert model, and for MPNet, it is set to 128 during the local training by participating clients.

Figures 11 and 12 depict the performance of MeanCache as FL training progresses. The x-axis represents the training round, while the y-axis shows the performance metrics such as F-score, precision, recall, and accuracy of the global model ( $W_{global}^{t+1}$ ). As illustrated in Figure 11, the F-score for MPNet increases from 0.82 to 0.88, and for Albert, it rises from 0.83 to 0.86, as shown in Figure 12. Similarly, precision for MPNet significantly increases from 0.74 to 0.85, as depicted in Figure 11, and for Albert, it increases from 0.74 to 0.81, as demonstrated in Figure 12. Given that MPNet is a more robust transformer architecture compared to Albert, it is also observed during our training that MPNet outperforms Albert, exhibiting superior learning in FL settings.

**Summary.** FL training increases 11% precision of MeanCache for MPNet and a 7% increase for Albert. The performance of the embedding model to generate high-quality embeddings can improve in a privacy-preserving manner using FL training.

#### F. Cosine Similarity Threshold Impact on Semantic Matching

Semantic matching for a new user query begins by generating the embeddings of the user’s query ( $E_q$ ) using the

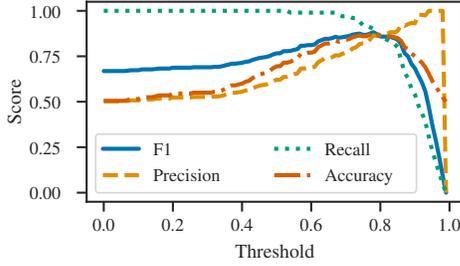


Fig. 14: MeanCache identifies an optimal threshold of 0.78 for Albert.

embedding model. The cosine similarity ( $\theta$ ) is computed with the cached embeddings. If the cosine similarity  $\theta$  exceeds the threshold  $\tau$ , the cache is hit and the response to the user query is returned from the cache. Therefore, the cosine similarity threshold  $\tau$  is crucial in determining the similarity between a user query and cached entries. A low threshold value of  $\tau$  can lead to false hits (incorrect matches), while a high threshold might overlook the appropriate matches (*i.e.*, false cache misses or false negatives).

To illustrate this, MeanCache varies the threshold  $\tau$  from 0 to 9 and evaluates the performance metrics F-score, precision, recall, and accuracy with an equal distribution of duplicate and non-duplicate queries from the validation data to avoid bias. Figures 13 and 14 show how the cosine similarity threshold ( $\tau$ ) affects MeanCache’s performance. The x-axis represents the threshold  $\tau$  values, and the y-axis denotes the performance metrics. For instance, at a 0.3 threshold, MeanCache’s semantic matching accuracy using MPNet is 57%, with a precision of 54% as shown in Figure 13. Similarly, with Albert at the same threshold, the accuracy is 55%, and the precision is 53% (Figure 14). Precision typically improves with an increase in threshold. However, beyond a certain point, increasing the threshold  $\tau$  leads to a decline in F score, accuracy, and recall.

For MPNet, the optimal threshold  $\tau$  is identified at 0.83, achieving an F1 score of 0.89, precision of 0.92, and accuracy of 0.90 (Figure 13). For Albert, the optimal threshold is 0.78, with an F1 score of 0.88.

**Summary.** GPTCache’s suggested threshold of 0.7 will result in suboptimal performance during semantic matching. The optimal threshold  $\tau$  values varies with the embedding model. MeanCache optimally adjusts the threshold based on user data, outperforming GPTCache’s suggested threshold by 16% in precision and 4% in F score for MPNet, and by 10% in precision and 2% in F score for Albert.

### G. Infeasibility of Embedding Generation with Llama 2

GPTCache [20] recommend using Llama for generating embeddings to enhance GPTCache’s performance. However, Llama 2’s embedding computation is expensive in terms of inference time, requires substantial storage, and incurs considerable overhead during semantic searches. For example,

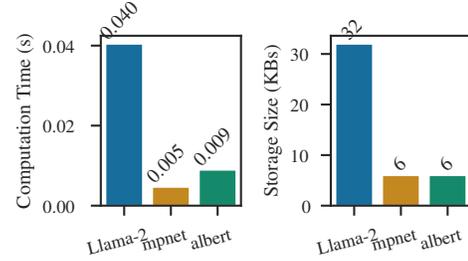


Fig. 15: Llama 2 takes significantly longer to compute embeddings and requires substantially more storage space than Albert and MPNet.

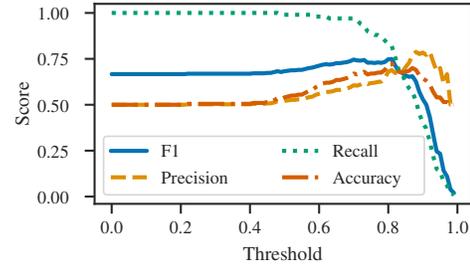


Fig. 16: Llama 2 does not perform well in semantic matching, even at optimal thresholds.

Llama 2 with its 7 billion parameters, demands 30 GB of memory [50], whereas Albert and MPNet require only 43 MB and 420 MB, respectively.

To highlight the impracticality to generate query embeddings with Llama 2, we compare the embedding computation time, embedding storage space requirement of Llama 2, Albert and MPNet transformer models. Figure 15 shows the average time to compute the embeddings of a single query and storage requirements for the embeddings. Figure 15 shows that the average embedding computation time of Llama 2 is 0.04 seconds, while for Albert and MPNet the average computation time is 0.005 and 0.009 seconds respectively. Single query embeddings generated from Llama 2 takes approximately 32 KBs of space and embeddings generated by both MPNet and Albert only take 6 KBs, as shown in Figure 15.

Furthermore, we evaluate the performance of the Llama 2 on embedding generation and semantic matching. Figure 16 shows that the performance of Llama 2 with different cosine similarity threshold ( $\tau$ ) and corresponding performance metric. We can see that the performance of Llama 2 is not good even with the optimal cosine similarity threshold, as also noted by researchers [51]. The maximum F1 score achieved by Llama 2 is 0.75 which is quite low when compared with the optimal thresholds scores from the Figures 13 and Figure 14.

Smaller models tailored for specific tasks often surpass larger models in efficiency [23], [24], [25]. Thus, diverging from GPTCache’s approach, we advocate for adopting smaller yet efficient embedding models for semantic caching. These models not only ensure optimal performance but also minimize

the semantic cache’s overhead on users, featuring lower inference demands and reduced output embedding sizes, thereby facilitating deployment on edge devices.

**Summary.** Llama 2 is not viable for generating embeddings. Future enhancements might improve its performance to generate embeddings, but the computational demands, semantic search duration, and storage requirements will likely remain elevated.

## V. RELATED WORK

Several caching systems are proposed to optimize the performance. Study [18] suggests a two-tier dynamic caching architecture for web search engines to enhance response times in hierarchical systems. Utilizing LRU eviction at both levels, they demonstrate how the second-tier cache can significantly lower disk traffic and boost throughput. Researchers in [52] propose a three-level index organization and [53] propose a three-tier caching. Another study [14] examined two real search engine datasets to explore query locality, aiming to develop a caching strategy based on this concept. Their analysis centered on query frequency and distribution, assessing the feasibility of caching at various levels, such as server, proxy, and client side. A novel caching technique called Probability Driven Cache (PDC) is proposed in [13] to optimize the performance of search engines by using the probability of query repetition to decide whether to cache the query. PDC uses the probability of a query to be repeated to decide whether to cache the query or not. A different approach is presented in [54], which proposes the Static Dynamic Cache (SDC) to exploit temporal and spatial locality present in the query stream, avoiding redundant processing and saving computational resources. Efficient caching designs for web search engines are explored in [17], where static and dynamic caching strategies are compared, weighing the benefits of caching query results against posting lists. Challenges in large-scale search engines, which process thousands of queries per second across vast document collections, are examined in [55], focusing on index compression, caching optimizations, and evaluating various inverted list compression algorithms alongside caching policies such as LRU and LFU.

All of these studies focus on caching systems designed for traditional search engines that process keyword queries (i.e., keyword matching) and return a list of links as a response. However, when applied to LLM-based web servers or APIs, these caching systems do not provide a single concise response and may yield many false results. Moreover, such caching techniques fail to capture the semantic similarity among repeated queries, leading to a significantly low hit rate. Server-side caching for services based on LLMs is proposed in [19] and [20], aiming to reduce the massive computational cost of LLMs. In particular, the approach in [19] checks if a new query is semantically similar to any existing queries in the cache. If a match is found, the cached response is returned; otherwise, a model multiplexer selects the most suitable LLM.

While these techniques can handle semantic similarity among queries and provide a single concise response, they raise privacy concerns as user queries are stored on external servers. Additionally, these techniques are static and unable to adapt to individual user behavior. Users may still be charged for the query, even if it is served from the cache. Therefore, a user-centered semantic cache that operates on the user side is needed, providing benefits in terms of privacy, cost, and latency. This cache should be able to detect semantic similarity among queries and adapt to each user’s behavior while preserving privacy. MeanCache offers these benefits without compromising user privacy.

## VI. CONCLUSION

MeanCache introduces the first user-centric semantic cache designed for LLM-based web services, such as ChatGPT. In MeanCache, clients collaboratively train a global embedding model using FL on their local data, ensuring user privacy. After aggregation, the global model produces high-quality embeddings for effective semantic matching. When a new query from the user matches a previous one, MeanCache semantically compares it with the user’s local cache and retrieves the most relevant results. This approach reduces the computational cost of LLM services, enhances bandwidth and latency, and conserves the user’s query quota. Even with compressed embeddings that save 83% of storage space, MeanCache outperforms existing baseline. With its distributed cache design, MeanCache offers a solution to reduce up to one-third of LLM query inference costs for semantically similar queries on the user side.

**Acknowledgment.** We thank anonymous reviewers for providing valuable and constructive feedback to help improve the quality of this work. This work was supported by Amazon - Virginia Tech Initiative in Efficient and Robust Machine Learning. We also thank the Advanced Research Computing Center at Virginia Tech and the Flower FL framework for their support in building and evaluating this work.

## REFERENCES

- [1] “Introducing chatgpt,” <https://openai.com/blog/chatgpt>, (Accessed on 01/18/2024).
- [2] “Google ai updates: Bard and new ai features in search,” <https://blog.google/technology/ai/bard-google-ai-search-updates/>, (Accessed on 01/18/2024).
- [3] “Introducing claude 2.1 \ anthropic,” <https://www.anthropic.com/news/claude-2-1>, (Accessed on 01/18/2024).
- [4] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.
- [5] “Learning human actions on computer applications,” <https://www.rabbit.tech/research>, (Accessed on 01/19/2024).
- [6] “Arc max is the popular browser’s new suite of ai tools - the verge,” <https://www.theverge.com/2023/10/3/23898907/arc-max-ai-browser-mac-ios>, (Accessed on 01/19/2024).
- [7] E. Frantar, S. Ashkboos, T. Hoefler, and D. Alistarh, “OPTQ: accurate quantization for generative pre-trained transformers,” in *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. [Online]. Available: <https://openreview.net/pdf?id=tcBBPnfwxS>

- [8] A. Tseng, J. Chee, Q. Sun, V. Kuleshov, and C. De Sa, “Quip#: Even better llm quantization with hadamard incoherence and lattice codebooks,” *arXiv preprint arXiv:2402.04396*, 2024.
- [9] N. Sachdeva, B. Coleman, W.-C. Kang, J. Ni, L. Hong, E. H. Chi, J. Caverlee, J. McAuley, and D. Z. Cheng, “How to train data-efficient llms,” *arXiv preprint arXiv:2402.09668*, 2024.
- [10] S. Ma, H. Wang, L. Ma, L. Wang, W. Wang, S. Huang, L. Dong, R. Wang, J. Xue, and F. Wei, “The era of 1-bit llms: All large language models are in 1.58 bits,” *arXiv preprint arXiv:2402.17764*, 2024.
- [11] “Perplexity pro,” <https://www.perplexity.ai/pro>, (Accessed on 03/01/2024).
- [12] “OpenAI Pricing,” <https://openai.com/pricing>, (Accessed on 01/19/2024).
- [13] R. Lempel and S. Moran, “Predictive caching and prefetching of query results in search engines,” in *Proceedings of the 12th international conference on World Wide Web*, 2003, pp. 19–28.
- [14] Y. Xie and D. O’Hallaron, “Locality in search engine queries and its implications for caching,” in *Proceedings. Twenty-First Annual Joint Conference of the IEEE Computer and Communications Societies*, vol. 3. IEEE, 2002, pp. 1238–1247.
- [15] E. P. Markatos, “On caching search engine query results,” *Computer Communications*, vol. 24, no. 2, pp. 137–143, 2001.
- [16] S. Podlipnig and L. Böszörményi, “A survey of web cache replacement strategies,” *ACM Computing Surveys (CSUR)*, vol. 35, no. 4, pp. 374–398, 2003.
- [17] R. Baeza-Yates, A. Gionis, F. Junqueira, V. Murdock, V. Plachouras, and F. Silvestri, “The impact of caching on search engines,” in *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, 2007, pp. 183–190.
- [18] P. C. Saraiva, E. Silva de Moura, N. Ziviani, W. Meira, R. Fonseca, and B. Ribeiro-Neto, “Rank-preserving two-level caching for scalable search engines,” in *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*, 2001, pp. 51–58.
- [19] B. Zhu, Y. Sheng, L. Zheng, C. Barrett, M. Jordan, and J. Jiao, “Towards optimal caching and model selection for large model inference,” in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. [Online]. Available: <https://openreview.net/forum?id=gd20oaZqqF>
- [20] F. Bang, “Gptcache: An open-source semantic cache for llm applications enabling faster answers and cost savings,” in *Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software (NLP-OSS 2023)*, 2023, pp. 212–218.
- [21] “Federated learning: Collaborative machine learning without centralized training data – google research blog,” <https://blog.research.google/2017/04/federated-learning-collaborative.html>, (Accessed on 04/01/2024).
- [22] K. Song, X. Tan, T. Qin, J. Lu, and T.-Y. Liu, “Mpnnet: Masked and permuted pre-training for language understanding,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 16 857–16 867, 2020.
- [23] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, “Training language models to follow instructions with human feedback,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 27 730–27 744, 2022.
- [24] G. Penedo, Q. Malartic, D. Hesslow, R. Cojocaru, H. Alobeidli, A. Cappelli, B. Pannier, E. Almazrouei, and J. Launay, “The refinedweb dataset for falcon llm: Outperforming curated corpora with web data only,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [25] Y. Du and L. Kaelbling, “Compositional generative modeling: A single model is not all you need,” *arXiv preprint arXiv:2402.01103*, 2024.
- [26] X. Wang, Q. Le, A. F. Khan, J. Ding, and A. Anwar, “Icl: An incentivized collaborative learning framework,” in *2024 IEEE International Conference on Big Data (BigData)*, 2024, pp. 94–103.
- [27] A. F. Khan, Y. Li, X. Wang, S. Haroon, H. Ali, Y. Cheng, A. R. Butt, and A. Anwar, “Towards cost-effective and resource-aware aggregation at edge for federated learning,” in *2023 IEEE International Conference on Big Data (BigData)*, 2023, pp. 690–699.
- [28] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [29] W. Gill, A. Anwar, and M. A. Gulzar, “TraceFL: Interpretability-Driven Debugging in Federated Learning via Neuron Provenance,” in *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. IEEE, 2025.
- [30] K. Pearson, “Liii. on lines and planes of closest fit to systems of points in space,” *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, vol. 2, no. 11, pp. 559–572, 1901. [Online]. Available: <https://doi.org/10.1080/14786440109462720>
- [31] H. Hotelling, “Analysis of a complex of statistical variables into principal components,” *Journal of educational psychology*, vol. 24, no. 6, p. 417, 1933.
- [32] I. T. Jolliffe and J. Cadima, “Principal component analysis: a review and recent developments,” *Philosophical transactions of the royal society A: Mathematical, Physical and Engineering Sciences*, vol. 374, no. 2065, p. 20150202, 2016.
- [33] “zilliztech/gptcache: Semantic cache for llms. fully integrated with langchain and llama\_index.” <https://github.com/zilliztech/gptcache>, (Accessed on 03/03/2024).
- [34] “Gptcache/examples/benchmark at main · zilliztech/gptcache,” <https://github.com/zilliztech/GPTCache/tree/main/examples/benchmark>, (Accessed on 03/04/2024).
- [35] D. J. Beutel, T. Topal, A. Mathur, X. Qiu, J. Fernandez-Marques, Y. Gao, L. Sani, H. L. Kwing, T. Parcollet, P. P. d. Gusmão, and N. D. Lane, “Flower: A friendly federated learning research framework,” *arXiv preprint arXiv:2007.14390*, 2020.
- [36] A. F. Khan, A. A. Khan, A. M. Abdelmoniem, S. Fountain, A. R. Butt, and A. Anwar, “Float: Federated learning optimizations with automated tuning,” in *Proceedings of the Nineteenth European Conference on Computer Systems*, ser. EuroSys ’24. New York, NY, USA: Association for Computing Machinery, 2024, p. 200–218. [Online]. Available: <https://doi.org/10.1145/3627703.3650081>
- [37] W. Gill, A. Anwar, and M. A. Gulzar, “FedDebug: Systematic Debugging for Federated Learning Applications,” in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, 2023, pp. 512–523.
- [38] —, “FedDefender: Backdoor Attack Defense in Federated Learning,” in *Proceedings of the 1st International Workshop on Dependability and Trustworthiness of Safety-Critical Systems with Machine Learned Components*, 2023, pp. 6–9.
- [39] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization in heterogeneous networks,” *Proceedings of Machine learning and systems*, vol. 2, pp. 429–450, 2020.
- [40] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, “Federated learning with matched averaging,” in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=BkluqSFDS>
- [41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [42] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence embeddings using Siamese BERT-networks,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, K. Inui, J. Jiang, V. Ng, and X. Wan, Eds. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 3982–3992. [Online]. Available: <https://aclanthology.org/D19-1410>
- [43] V. Karpukhin, B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W.-t. Yih, “Dense passage retrieval for open-domain question answering,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, B. Webber, T. Cohn, Y. He, and Y. Liu, Eds. Online: Association for Computational Linguistics, Nov. 2020, pp. 6769–6781. [Online]. Available: <https://aclanthology.org/2020.emnlp-main.550>
- [44] J. Ni, G. Hernandez Abrego, N. Constant, J. Ma, K. Hall, D. Cer, and Y. Yang, “Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models,” in *Findings of the Association for Computational Linguistics: ACL 2022*, S. Muresan, P. Nakov, and A. Villavicencio, Eds. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 1864–1874. [Online]. Available: <https://aclanthology.org/2022.findings-acl.146>
- [45] M. Henderson, R. Al-Rfou, B. Strope, Y.-H. Sung, L. Lukács, R. Guo, S. Kumar, B. Miklos, and R. Kurzweil, “Efficient natural language response suggestion for smart reply,” *arXiv preprint arXiv:1705.00652*, 2017.
- [46] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, “ALBERT: A lite BERT for self-supervised learning of language representations,” *CoRR*, vol. abs/1909.11942, 2019. [Online]. Available: <http://arxiv.org/abs/1909.11942>

- [47] “Diskcache: Disk backed cache — diskcache 5.6.1 documentation,” <https://grantjenks.com/docs/diskcache/>, (Accessed on 04/01/2024).
- [48] D. Avdiukhin and S. Kasiviswanathan, “Federated learning under arbitrary communication patterns,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 425–435.
- [49] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, “Tackling the objective inconsistency problem in heterogeneous federated optimization,” *Advances in neural information processing systems*, vol. 33, pp. 7611–7623, 2020.
- [50] “Sizing guide - nvidia docs,” <https://docs.nvidia.com/ai-enterprise/workflows-generative-ai/0.1.0/sizing-guide.html>, (Accessed on 01/25/2024).
- [51] “[user] embedding doesn’t seem to work? · issue #899 · ggerganov/llama.cpp,” <https://github.com/ggerganov/llama.cpp/issues/899>, (Accessed on 01/18/2024).
- [52] R. Baeza-Yates and F. Saint-Jean, “A three level search engine index based in query log distribution,” in *String Processing and Information Retrieval: 10th International Symposium, SPIRE 2003, Manaus, Brazil, October 8-10, 2003. Proceedings 10*. Springer, 2003, pp. 56–65.
- [53] X. Long and T. Suel, “Three-level caching for efficient query processing in large web search engines,” in *Proceedings of the 14th international conference on World Wide Web*, 2005, pp. 257–266.
- [54] T. Fagni, R. Perego, F. Silvestri, and S. Orlando, “Boosting the performance of web search engines: Caching and prefetching query results by exploiting historical usage data,” *ACM Transactions on Information Systems (TOIS)*, vol. 24, no. 1, pp. 51–78, 2006.
- [55] J. Zhang, X. Long, and T. Suel, “Performance of compressed inverted list caching in search engines,” in *Proceedings of the 17th international conference on World Wide Web*, 2008, pp. 387–396.