# Accelerating High-Dimensional Nearest Neighbor Search with Dynamic Query Preference

Yunjun Gao
Zhejiang University
gaoyj@zju.edu.cn

Ruijie Zhao
Zhejiang University
ruijie.zhao.work@gmail.com

Zhonggen Li
Zhejiang University
zgli@zju.edu.cn

Baihua Zheng
Singapore Management University
bhzheng@smu.edu.sg

Yifan Zhu
Zhejiang University
xtf_z@zju.edu.cn

Zhaoqiang Chen
Zhejiang University
chenzhaoqiang1@huawei.com

## ABSTRACT

Approximate Nearest Neighbor Search (ANNS) is a crucial operation in databases and artificial intelligence. Current graph-based ANNS methods, such as HNSW and NSG, have shown remarkable performance but are designed under the assumption of a uniform query distribution. However, in practical scenarios, user preferences and query temporal dynamics lead to some queries being searched for more frequently than others. To fully utilize these characteristics, we propose **DQF**, a novel **D**ual-Index **Q**uery **F**ramework. This framework comprises a dual-layer index structure and a dynamic search strategy based on a decision tree. The dual-layer index structure comprises a hot index for high-frequency nodes and a full index for the entire dataset, allowing for the separate management of hot and cold queries. Furthermore, we propose a dynamic search strategy that employs a decision tree to adapt to the specific characteristics of each query. The decision tree evaluates whether a query is of the high-frequency type to detect the opportunities for early termination on the dual-layer, avoiding unnecessary searches in the full index. Experimental results on four real-world datasets demonstrate that the Dual-Index Query Framework achieves a significant speedup of 2.0–5.7× over state-of-the-art algorithms while maintaining a 95% recall rate. Importantly, it does not require full index reconstruction when query distributions change, underscoring its efficiency and practicality in dynamic query distribution scenarios.

## 1 INTRODUCTION

Nearest Neighbor Search (NNS) in high-dimensional spaces plays a crucial role in a wide range of applications, including information retrieval [9, 23, 39], recommendation systems [6, 29, 30], and retrieval-augmented generation [10, 38]. However, due to the well-known curse of dimensionality [20, 34], exact nearest neighbor search becomes computationally expensive and inefficient, especially at large scales. To address this issue, numerous approximate nearest neighbor search (ANNS) methods have been proposed, including hash-based [18, 19, 25], quantization-based [22, 28, 33], tree-based [3, 4, 11], and graph-based methods [32, 37].

Among these, graph-based methods have demonstrated particularly strong performance, offering fast search speeds while maintaining high recall. These methods typically construct an index in the form of a graph, where nodes represent data points and edges capture proximity relationships. The search typically starts from an initial node and iteratively expands to neighboring nodes, navigating the graph toward the nearest neighbors [26]. Recent research has focused on improving search performance by optimizing edge selection and pruning strategies [13, 14, 21, 27]. For example, HNSW [27] constructs a hierarchical, multi-layered graph to speed up the search by starting at the top layer and descending to lower layers. NSG [14], on the other hand, utilizes the Relative Neighborhood Graph (RNG) property [31] to ensure that each search step moves closer to the query point. Both methods aim to refine the graph structure and edge selection strategies for better efficiency.

Despite their effectiveness, these methods often assume that all data points have equal query probability, overlooking critical real-world factors such as user preferences and temporal dynamics. In practice, user queries are far from uniform. For example, YouTube videos with more plays are more likely to be watched [16], and frequently visited websites are more likely to appear in Google search results [24]. These behaviors often follow Zipf's law [2], where the frequency of an item is inversely proportional to its rank. As shown in Figure 1, the query frequency of data points decreases with increasing ranks. Therefore, optimizing for user query preference, especially for high-frequency or trending queries, is crucial for improving retrieval efficiency and user satisfaction.

A straightforward approach to incorporate temporal query dynamics is to factor them into the graph construction process. For instance, PANNS [36] integrates recency by assigning higher weights to recently accessed nodes, increasing their likelihood of being connected in the graph, while pruning long-idle nodes. However,
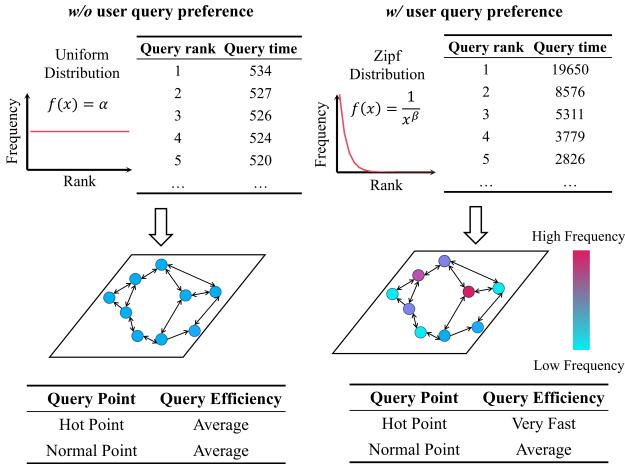
**Figure 1: Query Distribution Patterns: Uniform vs. Zipf.**

this strategy has two key limitations. First, when an index becomes outdated due to temporal query dynamics, the entire index has to be reconstructed. This process is resource-intensive, time-consuming, and inefficient in dynamic environments where query patterns evolve continuously, particularly for large-scale datasets. Second, this strategy considers only query recency, neglecting user query preferences. In real-world applications, users often search for current trending topics and exhibit a clear preference for high-frequency data. Therefore, a more comprehensive approach is needed - one that incorporates both the query timeliness and user preference into the indexing and search processes. Achieving this, however, raises three primary challenges:

*Challenge I: How can we effectively leverage user query preferences?* A major challenge lies in efficiently handling user query preferences in vector data. Unlike non-vectorized key-value data, where caching mechanisms can be employed to optimize access, vector-based search requires similarity computations that are far more complect. Consequently, exact matches and traditional caching techniques are infeasible. To tackle this challenge, we propose a novel dual-layer index structure. The upper layer, referred to as the hot index, contains high-frequency nodes that are queried most often. The lower layer, known as the full index, holds the entire dataset, including those in hot index. By prioritizing the hot index during search, we can rapidly serve the most likely queries. This structure separates hot and cold queries, allowing computational resources to be allocated more effectively. As a result, frequently queried data can be retrieved more quickly, improving responsiveness without sacrificing accuracy.

*Challenge II: How to adapt to dynamic query distributions?* In real-world applications, query distributions are constantly evolving, demanding an indexing method that can dynamically adapt to these changes. Existing approaches like PANNS require full index reconstruction to stay current, which is prohibitively slow, taking days on billion-scale dataset [21]. In contrast, our framework rebuilds only the upper-layer index. Since the upper-layer index is substantially smaller in scale compared to the complete index, our method achieves a speed-up of hundreds of times over traditional

methods. This efficiency allows us to update the index promptly, ensuring it reflects the latest query trends. Consequently, our system can maintain agility and responsiveness, keeping up with the dynamic nature of query distributions and providing timely and relevant search results.

*Challenge III: How can we search more efficiently?* A key challenge in enhancing search efficiency is that existing search methods don't distinguish between high-frequency and low-frequency nodes. Existing methods treat all data points equally during search, failing to distinguish between high- and low-frequency queries. This results in unnecessary overhead, especially after searching the hot index. Furthermore, traditional early termination methods provide minimal grains (~5%) in uniform query scenarios [7]. We introduce a decision tree-based dynamic routing strategy that distinguishes between high- and low-frequency queries based on observed search patterns. For instance, high-frequency queries typically yield stable top-$k$ results from the hot index, while low-frequency queries require deeper traversal. By exploiting these differences, our method achieves 2.0-5.7× speedups through dynamic search adjustments that reduce overhead and improve performance for both query types.

In summary, this work makes the following contributions:

1) We propose a novel **Dual-Index Query Framework** that integrates user query references for ANNS optimization.
2) We design a two-layered index structure to leverage user query preferences. The upper layer stores high-frequency nodes, while the lower layer maintains the complete dataset. This strategy allows us to manage hot and cold queries separately and prioritize high-frequency queries during the search process.
3) Our method rebuilds only the upper-layer index to adapt to shifting query distributions, achieving significant speedups over full-index reconstruction.
4) We introduce a decision tree-based query strategy that dynamically adjusts the search process, optimizing efficiency and accuracy across varying query frequencies.
5) Experiments on four real-world datasets show that our Dual-Index Query Framework maintains a 95% recall rate and achieves a 2.0–5.7× speedups over the state-of-the-art methods such as NSSG algorithm. Notably, it does not require full index reconstruction under changing query distributions.

## 2 RELATED WORK

In this section, we review prior work on approximate nearest neighbor search (ANNS) and indexing methods with temporal adaptation.

### 2.1 Non-Graph-Based Methods

Non-graph-based methods have been widely studied and applied in the field of information retrieval. Hash-based methods, such as Locality-Sensitive Hashing (LSH) [17], use hash functions to map high-dimensional data into low-dimensional hash codes, enabling fast similarity searches through bitwise operations. Quantization-based methods, like Product Quantization (PQ) [22] and Optimized Product Quantization (OPQ) [15], compress vectors into quantized

codes to reduce storage and computational requirements. Tree-based methods, including KDTree [4] and VPTree [11], partition the data space into hierarchical tree structures for efficient search. However, these methods face significant challenges in high-dimensional spaces. For example, tree-based methods suffer from the "curse of dimensionality", which diminishes the effectiveness of hierarchical partitioning as the dimensionality increases. Similarly, hash- and quantization-based methods often yield lower result quality in high-dimensional spaces due to hash collisions and quantization errors, respectively. These limitations have motivated the development of more sophisticated approaches, with graph-based methods emerging as a promising direction.

## 2.2 Graph-Based Methods

Graph-based methods have gained significant attention due to their effectiveness in handling high-dimensional and large-scale datasets. These methods construct a graph where nodes represent data points and edges represent similarities between them. NSW [26] builds a navigable small-world graph by connecting each new node to its nearest neighbors. HNSW [27] improves upon NSW by introducing a hierarchical structure, enabling faster search with logarithmic complexity. NSG [14] builds a sparse graph using a pruning strategy based on the Monotonic Relative Neighborhood Graph (MRNG) theory. NSSG [13] further optimizes the graph construction process with a satellite system graph (SSG) pruning strategy, which ensures a more even distribution of out-edges and adaptively adjusts graph sparsity. These graph-based methods excel in search speed and recall, particularly on large-scale, high-dimensional datasets.

However, more existing graph-based approaches assume uniform query distributions, which limits their effectiveness in real-world scenarios involving user preferences or temporal shifts in query patterns. This rigidity hinders their adaptability to dynamic environments where query behavior frequently evolves.

## 2.3 Temporal Adaptation in Indexing

The timeliness of queries presents unique challenges for ANNS. PANNS [36] addresses this by integrating temporal information into graph construction. It emphasizes the recency of data points, increasing the likelihood that recently accessed nodes are connected in the graph. It also introduces a fully parameterized beam search algorithm to optimize search performance. While PANNS demonstrates improved responsiveness to recent queries, it has key limitations. Most notably, it struggles to adapt dynamically to evolving query patterns. When the query distribution changes, the index often requires full reconstruction, a process that is both resource-intensive and time-consuming, especially on large-scale datasets. This makes it impractical for real-time or near-real-time systems.

The dynamic nature of real-world data necessitates efficient mechanisms for incorporating emerging trends without rebuilding the entire index. Frequent index updates are computationally expensive and, if not handled efficiently, can lead to suboptimal search performance. Therefore, developing methods that can adapt to shirting query patterns on the fly, without the overhead of full index reconstruction, remains a critical and open challenge.

## 3 BACKGROUND

In this section, we introduce the background of approximate nearest neighbor search and Zipf's law.

## 3.1 Approximate Nearest Neighbor Search

Approximate Nearest Neighbor Search (ANNS) aims to efficiently retrieve points that are close to a query point in a high-dimensional space, trading off a small amount of accuracy for significant gains in speed and scalability. Formally, given a dataset $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n\}$ where each $\mathbf{x}_i \in \mathbb{R}^d$, and a query point $\mathbf{q} \in \mathbb{R}^d$, exact nearest neighbor search (NNS) is to find a set $\mathcal{N}_k(\mathbf{q})$ of $k$ points from $\mathcal{D}$ that are closest to $\mathbf{q}$ under a certain distance metric $\text{dist}(\cdot, \cdot)$. Mathematically, NNS can be defined as:

$$\mathcal{N}_k(\mathbf{q}) = \arg \min_{\mathcal{S} \subseteq \mathcal{D}, |\mathcal{S}|=k} \sum_{\mathbf{x} \in \mathcal{S}} \text{dist}(\mathbf{q}, \mathbf{x}). \tag{1}$$

However, due to the computational complexity of exact NNS, especially in high-dimensional spaces, approximate methods are often employed. In this paper, we focus on the approximate $k$-nearest neighbor search, where the goal is to return a set $\mathcal{A}_k(\mathbf{q})$ such that:

$$\mathcal{A}_k(\mathbf{q}) \subseteq \mathcal{D}, \ \max_{\mathbf{x} \in \mathcal{A}_k(\mathbf{q})} \text{dist}(\mathbf{q}, \mathbf{x}) \leq (1 + \epsilon) \min_{\mathbf{x} \in \mathcal{D}} \text{dist}(\mathbf{q}, \mathbf{x}), \tag{2}$$

for a given approximation factor $\epsilon > 0$.

The performance of ANNS algorithms is typically evaluated using the recall metric, which measures the proportion of true nearest neighbors retrieved by the algorithm. Specifically, the recall@$k$ is defined as:

$$\text{Recall@}k = \frac{|\mathcal{A}_k(\mathbf{q}) \cap \mathcal{N}_k(\mathbf{q})|}{k}. \tag{3}$$

In this paper, we evaluate the performance of different ANNS algorithms based on both recall and computational efficiency.

## 3.2 Applications of Zipf's Law

Zipf's law originates from the observation of word frequency distributions in natural language corpora and has found extensive applications in information retrieval. The law states that in a large set of items, the frequency of any item is inversely proportional to its rank in the frequency table. This relationship can be empirically expressed as:

$$f(r) \sim r^{-\beta} \tag{4}$$

where $r$ denotes the rank of a word's frequency, $\beta$ is a constant that adjusts the Zipf distribution, and $f(r)$ represents the frequency of occurrence.

The prevalence of Zipf's law extends beyond corpora. Aaron Clauset et al. [8] rigorously analyzed 24 real-world datasets and confirmed the widespread occurrence of Zipf's distribution in practical scenarios. This pattern significantly impacts the field of querying. Research by Phillippa Gill et al. [16] on YouTube video viewership uncovered a phenomenon of concentrated queries on popular videos. Fabrizio Lillo et al. [24] discovered that the distribution of Google search volumes adheres to Zipf's law, indicating a clear preference in users' information-seeking behavior. Furthermore, Lada A. Adamic et al. [1] provided additional evidence that web search queries also follow this pattern. In the context of information

retrieval and recommendation systems, Zipf's law is evident in the distribution of user queries.

These empirical observations reveal a critical mismatch between real-world query distributions and the uniform-access assumption underlying most existing ANNS indexes. Under Zipf's workloads, a small fraction of data points (the head of the distribution) receives the overwhelming majority of queries, while the vast tail is rarely accessed. Traditional methods that treat all points equally waste both memory and computation on rarely queried items and fail to prioritize fast access to high-demand ones. Recognizing this, we propose to physically separate the hot, high-frequency items from the cold, low-frequency ones within the index structure. Our framework maintains a compact, high-performance layer for Zipf's head and a comprehensive but secondary layer for the long tail. This dual-layer structure allows the system to allocate resources based on the actual query frequency, improving both efficiency and responsiveness. Moreover, it enables rapid adaptation to evolving query patterns, a crucial requirement in dynamic real-world environments. This insight serves as the foundation for our proposed dual-layer framework.

## 4 METHODOLOGY

In this section, we provide an overview of our proposed Dual-Index Query Framework, which is designed to enhance the efficiency and effectiveness of high-dimensional nearest neighbor queries by leveraging user query preferences and temporal dynamics. The framework consists of two main components: a dual-layer index structure and a dynamic search strategy powered by a decision tree. The overall architecture of the framework is illustrated in Figure 2.

### 4.1 Framework Overview

To distinguish between high-frequency and low-frequency nodes, the dual-layer index structure comprises a hot index and a full index. The hot index stores high-frequency nodes that are frequently accessed by users, while the full index maintains a complete dataset. This separation allows the index to manage hot and cold queries separately. During the search process, the framework prioritizes the hot index to quickly retrieve content that aligns with current user interests, thereby improving search speed and resource utilization.

However, due to the lack of recall guidance, we cannot determine whether we need to continue searching in the full index after completing the search of the hot index. Therefore, the dynamic search strategy utilizes a decision tree to recognize the characteristics of each query. When the search in the hot index is completed, the framework employs the decision tree to determine whether to continue searching in the full index or to terminate the search early. This decision is based on various features extracted during the search process, such as distance counts and update counts. By dynamically adjusting the search strategy, the framework efficiently allocates computational resources, reducing unnecessary computations for high-frequency queries while ensuring thorough results for low-frequency queries.

The framework also addresses the challenge of dynamic query distributions by focusing on rebuilding the upper-layer index (hot index). Since the hot index is much smaller than the complete index, this approach significantly reduces the time and resources required for index reconstruction, enabling the index to promptly reflect the latest query trends.

In summary, our Dual-Index Query Framework provides a comprehensive solution to the challenges of leveraging user query preferences and handling dynamic query distributions in high-dimensional nearest neighbor search. It combines the efficiency of a dual-layer index structure with the adaptability of a dynamic search strategy to deliver improved performance in dynamic and preference-driven query scenarios.

### 4.2 Dual-layer Index Construction

The dual-layer index structure is designed to address the challenges of efficiently handling both high-frequency and low-frequency queries in high-dimensional spaces. It separates the hot index, which focuses on high-frequency nodes, from the full index, which covers low-frequency queries. This separation allows the index to prioritize high-frequency queries, improving search speed for popular content and boosting overall performance.

Another key innovation of this index is its ability to dynamically adapt to changing query patterns and user preferences without requiring full-scale index reconstruction. This is particularly important in dynamic environments where query distributions can shift rapidly, making it impractical to rely on static indexes. The hot index can be quickly updated and rebuilt to reflect the latest query trends. Since the hot index is much smaller than the full index, this process is rapid and resource-efficient, ensuring the index maintains high-quality search results even as user interests evolve.

*4.2.1 Full Index.* The full index serves as the foundation of our dual-layer index structure, offering comprehensive coverage of the entire dataset. It is constructed using the Satellite System Graph (SSG) [13] as a baseline. SSG is a graph structure where each node's outgoing edges are evenly distributed in all directions, ensuring efficient navigation and search capabilities. When the search process moves from the hot index to the full index, the full index provides a robust and complete foundation for continuing the search process.

A key aspect of SSG is its pruning strategy, which ensures the even distribution of edges while maintaining navigability. This strategy removes redundant edges and retains only those that contribute to the graph's navigability and efficiency.

The algorithm for the SSG pruning strategy is shown as Algorithm 1. In this algorithm, we utilize EFANNA [12] to quickly construct a pre-built k-nearest neighbor graph (KNNG). EFANNA employs a divide-and-conquer strategy along with NN-descent optimizations to efficiently build the KNNG. For each node $p$, the candidate set $C$ is formed by including not only the direct neighbors from the KNNG but also the neighbors of those neighbors (lines 3-8). This candidate set is then sorted by distance (line 9). For any two edges $pq$ and $pr$, if the angle between them is less than the specified threshold $\alpha$, the longer edge is discarded (lines 10-20). This process ensures that the edges around each node are evenly distributed, maintaining near-minimal neighborhood coverage. This pruning strategy is crucial for the efficient construction of the full index as it reduces the graph's complexity while preserving its navigability.

By utilizing NSSG as the baseline for the full index, our framework ensures that the full index can effectively support the search process for low-frequency queries. The full index acts as a reliable
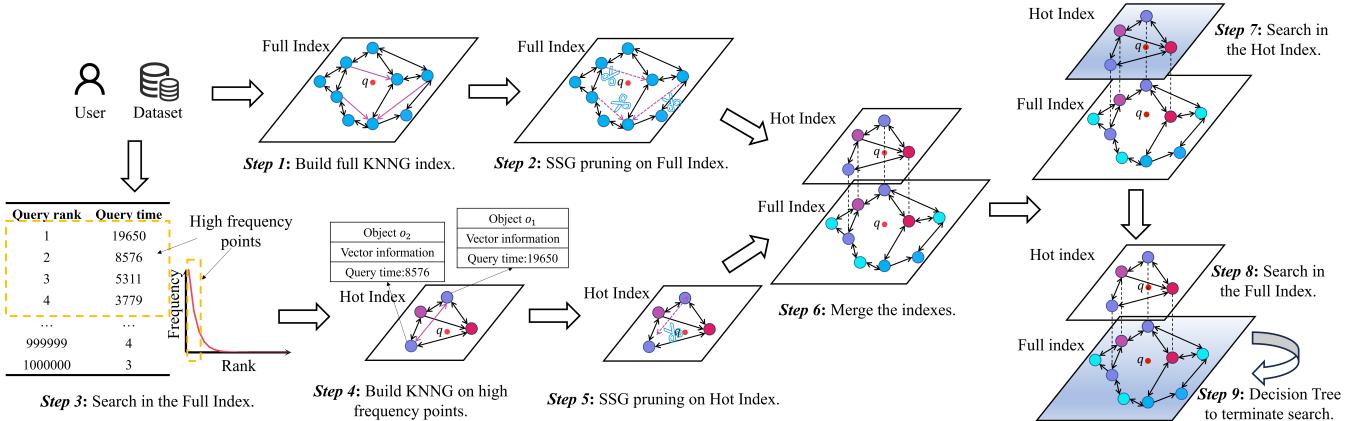
**Figure 2: The structure of our Dual-Index Query Framework, DQF**

---

**Algorithm 1** SSG Pruning Strategy

---

**Require:** pre-build KNNG $G_{knng}$, angle threshold $\alpha$
**Ensure:** Set of selected neighbors $P$
1:  $P \leftarrow \emptyset$
2:  $C \leftarrow \emptyset$
3:  **for all** neighbor $d_j$ of node $p$ in $G_{knng}$ **do**
4:      $C \leftarrow C \cup \{d_j\}$
5:      **for all** neighbor $d_k$ of node $d_j$ in $G_{knng}$ **do**
6:          $C \leftarrow C \cup \{d_k\}$
7:      **end for**
8:  **end for**
9:  Sort $C$ by distance in ascending order
10: **for all** node $d_j$ in sorted $C$ **do**
11:     flag $\leftarrow$ True
12:     **for all** node $d_k$ in $P$ **do**
13:         **if** $\cos \angle d_j p d_k > \cos \alpha$ **then**
14:             flag $\leftarrow$ False
15:         **end if**
16:     **end for**
17:     **if** flag = True **then**
18:         $P \leftarrow P \cup \{d_j\}$
19:     **end if**
20: **end for**

---

**Algorithm 2** Hot Index Construction

---

**Require:** Dataset $D$, Query trigger parameter $n\_query$, Hot index size $n\_idx$
**Ensure:** Hot index graph $G_{hot}$
1:  $query\_cnt \leftarrow 0$
2:  Use NSSG algorithm to build the full index graph $G_{full}$
3:  **while** True **do**
4:      Monitor and increment $query\_cnt$ for each node access
5:      **if** $query\_cnt > n\_query$ **then**    ▷ Trigger threshold for index update
6:          Sort nodes by query frequency in descending order
7:          Select top $n\_idx$ high-frequency nodes
8:          Construct $G_{hot}$ using NSSG algorithm on selected nodes
9:          Update $G_{hot}$ and replace the existing hot index
10:         $query\_cnt \leftarrow 0$
11:     **end if**
12: **end while**

---

backup, guaranteeing the completeness of search results and ensuring that all data points are accessible even if they are not part of the hot index.

*4.2.2 Hot Index.* Creating and managing the hot index presents unique challenges. Unlike the full index, the hot index must dynamically adapt to changing query patterns and user preferences. It needs to efficiently track and prioritize high-frequency nodes that are frequently accessed by users. Additionally, the hot index must balance the need for rapid updates with maintaining a structure that allows for efficient search operations.

To address these challenges, we have implemented a dynamic updating mechanism that triggers updates based on the query frequency. This strategy ensures that the hot index remains relevant

and reflective of current query trends. By maintaining a compact hot index, we not only reduce the complexity of the index but also enhance the search efficiency.

The hot index is designed to efficiently manage high-frequency queries by focusing on a subset of frequently accessed nodes within the dataset. It is constructed using a similar NSSG-based approach as the full index but is optimized to prioritize nodes that are more likely to be queried based on user preferences and temporal dynamics.

Algorithm 2 illustrates the process of constructing the hot index. In this algorithm, the hot index construction begins with initializing a query counter ($query\_cnt$) and building the full index graph ($G_{full}$) using the NSSG (lines 1-2). The main loop starts at line 3, where the index continuously monitors and increments the query counter for each node access. When the query counter exceeds the predefined threshold ($n\_query$) (line 5), the algorithm triggers an update process. This involves sorting the nodes by their query frequency in descending order (line 6) and selecting the top $n\_idx$ high-frequency nodes (line 7). The hot index graph ($G_{hot}$) is then reconstructed using the NSSG algorithm on these selected nodes

**Algorithm 3** Traditional Beam Search

**Require:** Graph index $G$, Query point $q$, Number of nearest neighbors $k$, Candidate pool size $l$
**Ensure:** Search results $Res$
 1: $L \leftarrow$ eps           ▷ Candidate pool initialized with entry points
 2: $V \leftarrow \emptyset$                           ▷ Set of visited nodes
 3: **while** $L \setminus V \neq \emptyset$ **do**
 4:      $p \leftarrow$ first unvisited node in $L$
 5:      $L \leftarrow L \cup \text{Neighbors}(p)$
 6:      $V \leftarrow V \cup \{p\}$
 7:      **if** $|L| > l$ **then**
 8:          Trim $L$ to retain only the $l$ closest nodes to $q$
 9:      **end if**
10: **end while**
11: **return** $k$ closest nodes from $L$ to $q$



**Figure 3: Comparison with NSSG, Traditional Beam Search, and Dynamic Search Using Decision Tree**

(line 8). Finally, the existing hot index is updated and replaced with the new $G_{hot}$, and the query counter is reset (lines 9-10). This dynamic updating mechanism allows the hot index to adapt efficiently to changing query patterns, providing more efficient access to high-frequency data points and enhancing the overall performance of the search process.

To summarize, the full index and hot index work together to provide a comprehensive and efficient search solution. The full index ensures comprehensive coverage and reliable backup for all data points, supporting both high-frequency and low-frequency queries. Meanwhile, the hot index focuses on frequently accessed nodes, prioritizing them to speed up the search for popular content. This dual-layer structure not only improves search speed and resource utilization but also ensures the index can adapt dynamically to changing query demands.

## 4.3 Dynamic Search Using Decision Tree

In this section, we introduce our approach to accelerating the dual-layer index search using a Decision Tree. Before diving into that, we first discuss beam search and explain why directly applying traditional beam search isn't effective for our dual-layer index.

*4.3.1 Traditional Beam Search.* The traditional beam search algorithm is a heuristic search strategy widely used in graph-based approximate nearest neighbor search. It maintains a queue of candidate nodes, iteratively expanding the most promising nodes while limiting the number of candidates to a fixed size.

Algorithm 3 demonstrates the process of the traditional beam search algorithm. The candidate pool $L$ is initialized with entry points (line 1), and the set of visited nodes $V$ is initialized as an empty set (line 2). The algorithm then enters a loop where it continues to expand the most promising nodes until all nodes in the candidate pool have been visited (line 3). In each iteration, the algorithm selects the first unvisited node $p$ from the candidate pool (line 4), expands its neighboring nodes and adds them to the candidate pool (line 5), and marks $p$ as visited (line 6). To maintain the size of the candidate pool within the specified limit $l$, the algorithm trims the pool to retain only the $l$ closest nodes to the query point $q$ if its size exceeds $l$ (lines 7-9). Finally, after all nodes in the candidate
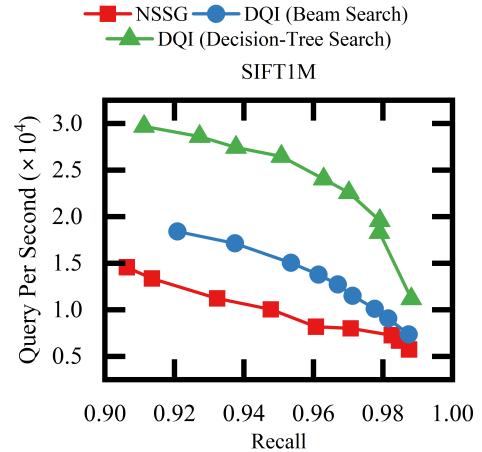
pool have been visited, the algorithm returns the $k$ closest nodes to $q$ from the candidate pool (line 10).

While the beam search algorithm can rapidly identify most target nodes, it does not fully exploit the timeliness and preference of queries. Due to the absence of clear recall guidance, after the algorithm completes the Hot Index search, it continues to search the Full Index until all nodes in the queue have been visited. This process is time-consuming and yields limited improvements in the recall rate for high-frequency nodes. As illustrated in Figure 3, DQI with traditional beam search results in only modest performance improvements over NSSG. Conversely, our dynamic decision tree search algorithm demonstrates substantial performance gains compared to NSSG. This highlights the limitations of applying traditional beam search algorithms directly to our Dual-layer Index, as they fail to fully utilize the advantages in handling high-frequency query nodes.

*4.3.2 Dynamic Search with Decision Tree.* While examining the limitations of traditional beam search, we identified a key issue: the inability to distinguish between high-frequency and low-frequency nodes. To address this, we explored enhancing the search strategy with a mechanism that could dynamically adapt to query characteristics. Our solution was to integrate a decision tree into the search process, enabling adaptive termination based on query-specific features. This approach not only distinguishes between high-frequency and low-frequency nodes but also optimizes resource allocation by reducing unnecessary computations.

Our dynamic search strategy enhances the traditional beam search by incorporating a decision tree to adaptively terminate the search process based on query-specific features. This reduces unnecessary computations and improves efficiency. The algorithm is shown as Algorithm 4.

This dynamic search approach first initializes the result list $L$ with entry points from $G_{hot}$ and marks these nodes as unvisited (lines 1-2). It then searches within $G_{hot}$, expanding candidate nodes and keeping the result list size within $s\_l$ (lines 3-10). After the
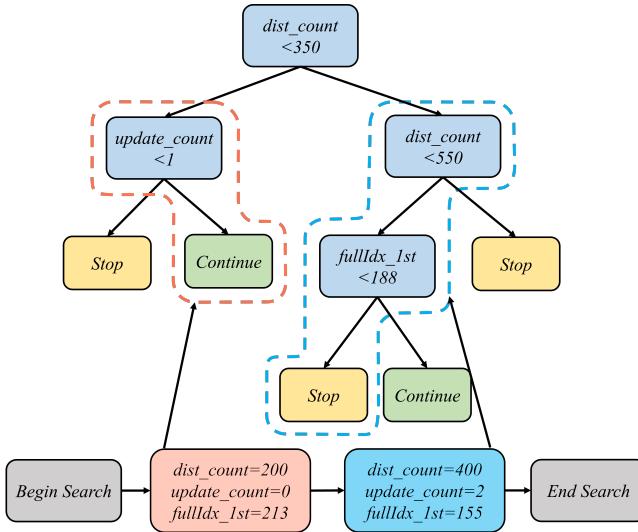
**Figure 4: Four-Layer Decision Tree for SIFT1M Dataset**

---

**Algorithm 4** Dynamic Search with Decision Tree

---

**Require:** Hot index $G_{hot}$, Full index $G_{full}$, Query point $q$, Number of nearest neighbors $k$, Hot index candidate pool size $s\_l$, Full index candidate pool size $l$, Decision tree judgment frequency $freq$

**Ensure:** Search results $Res$

1:  $L \leftarrow$ entry points of $G_{hot}$        ▷ Initialize result list
2:  $V \leftarrow \emptyset$        ▷ Initialize visited nodes set
3:  **while** $L \setminus V \neq \emptyset$ **do**
4:      $p \leftarrow$ first unvisited node in $L$
5:      $L \leftarrow L \cup$ Neighbors$(p)$ in $G_{hot}$
6:      $V \leftarrow V \cup \{p\}$
7:      **if** $|L| > s\_l$ **then**
8:         Trim $L$ to retain only the $s\_l$ closest nodes to $q$
9:      **end if**
10: **end while**
11: Reset visit status of nodes in $L$
12: $dist\_cnt \leftarrow 0$
13: **while** $L \setminus V \neq \emptyset$ **do**
14:     $p \leftarrow$ first unvisited node in $L$
15:     $L \leftarrow L \cup$ Neighbors$(p)$ in $G_{full}$
16:     $V \leftarrow V \cup \{p\}$
17:     $dist\_cnt \leftarrow dist\_cnt + 1$
18:     **if** $dist\_cnt$ mod $freq == 0$ **then**
19:        Use decision tree to decide whether to terminate the search
20:     **end if**
21:     **if** $|L| > l$ **then**
22:        Trim $L$ to retain only the $l$ closest nodes to $q$
23:     **end if**
24: **end while**
25: **return** $k$ closest nodes from $L$ to $q$

---

**Table 1: Decision Tree Features**

| Feature | Description |
|---|---|
| hotIdx_1st | The distance from the query point to the nearest node in the hot index. |
| hotIdx_1st_div_kth | The ratio of the first nearest node distance to the k-th nearest node distance in the hot index. |
| fullIdx_1st | The distance from the query point to the nearest node in the full index. |
| fullIdx_1st_div_kth | The ratio of the first nearest node distance to the k-th nearest node distance in the full index. |
| dist_count | The number of distance calculations performed during the search process. |
| update_count | The number of node accesses and updates during the search process. |

**Table 2: Feature Importance Across Datasets**

| Feature | SIFT1M | GIST1M | Crawl | Glove |
|---|---|---|---|---|
| hotIdx_1st | 11.3% | 9.1% | 10.6% | 11.4% |
| hotIdx_1st_div_kth | 7.9% | 8.2% | 11.9% | 9.0% |
| fullIdx_1st | 13.9% | 13.0% | 16.4% | 22.3% |
| fullIdx_1st_div_kth | 7.1% | 7.5% | 10.2% | 10.1% |
| dist_count | 46.3% | 43.7% | 39.0% | 35.0% |
| update_count | 13.5% | 18.5% | 12.0% | 12.2% |

$G_{hot}$ search, the algorithm resets the visited status of nodes in $L$ and sets up a distance counter $dist\_cnt$ (lines 11-12). Next, it continues the search in $G_{full}$, periodically using a decision tree to decide if the search should terminate early, thus cutting down on unnecessary computations (lines 13-24). This strategy is more efficient for high-frequency queries as it adjusts dynamically based on query characteristics, reducing waste of computing resources on low-frequency nodes.

In the dynamic search strategy, the decision tree uses six key features to distinguish between high-frequency and low-frequency query nodes, which can be categorized as follows: *(a) Distance-related features in the hot index:* hotIdx_1st represents the distance from the query point to the nearest node in the hot index, while hotIdx_1st_div_kth denotes the ratio of the first nearest node distance to the k-th nearest node distance within the hot index. Once the hot index search is completed, these distance-related features are determined. *(b) Distance-related features in the full index:* fullIdx_1st is the distance from the query point to the nearest node in the full index, and fullIdx_1st_div_kth is the corresponding ratio of the first nearest node distance to the k-th nearest node distance in the full index. *(c) Count-related features:* dist_count indicates the number of distance calculations performed during the search process, and update_count represents the number of node accesses and updates during the search. These features are detailed in Table 1. Additionally, Figure 4 illustrates a four-layer decision tree for the SIFT1M dataset as an example.

As shown in Table 2, the importance of each feature varies across datasets. Notably, the dist_count feature, which reflects the number of distance calculations during the search process, generally contributes significantly to the decision-making process. However, other features also play crucial roles in distinguishing between high-frequency and low-frequency query nodes. These features collectively help the decision tree adaptively refine the search strategy, thereby improving the overall efficiency and effectiveness of the dynamic search process.

To train the decision tree, we randomly sample historical queries, remove duplicates, and use the remaining queries in the full index phase. This ensures that the decision tree learns from a diverse set of queries representative of actual search patterns. During the search process, if the decision tree predicts that a query will likely receive future updates, it continues the search. Conversely, if no further updates are expected, the search process is terminated early. This approach optimizes the search process, ensuring that computational resources are allocated effectively. By leveraging these features and training methodology, the decision tree is enabled to dynamically adjust the search strategy to enhance the search efficiency and system responsiveness.

## 4.4 Complexity Analysis

In this section, we analyze the time complexity of the proposed Dual-Index Query Framework and determine the optimal Index Ratio (IR) that minimizes the overall complexity. Here, IR is defined as the ratio of the number of points in the hot index to the number of points in the full index. Traditional graph-based methods like NSSG have a time complexity of $O(n^{1/d} \log n)$ for search operations, where $d$ is the dimensionality of the data. This complexity arises due to the need to traverse multiple layers of the graph structure in high-dimensional spaces. However, in high-dimensional spaces, the term $n^{1/d}$ can be approximated as a constant factor when $d$ is sufficiently large, effectively simplifying the complexity to $O(\log n)$. This simplification is reasonable in many practical scenarios where the dimensionality $d$ is large enough to make $n^{1/d}$ close to 1.

Our framework builds upon this foundation and further optimizes the complexity by considering dynamic query preferences and user query patterns. The time complexity of the search process in our framework can be expressed as:

$$C(\text{IR}) = \log(\text{IR} \cdot n) + p \cdot \log n \qquad (5)$$

Here, IR represents the ratio of the hot index size to the full index size, $n$ is the total number of data points, and $p$ is the probability that a query cannot be resolved within the hot index alone. The first term $\log(\text{IR} \cdot n)$ corresponds to the complexity of searching within the hot index, while the second term $p \cdot \log n$ accounts for the complexity of searching within the full index.

The probability $p$ is derived from Zipf's distribution of query frequencies. Specifically, $p$ is calculated as:

$$p = 1 - \frac{\sum_{i=1}^{\text{IR} \cdot n} \frac{1}{i^\beta}}{\sum_{i=1}^{n} \frac{1}{i^\beta}} \qquad (6)$$

where $\beta$ is a parameter that adjusts Zipf's distribution. By approximating the summations using integrals, we can simplify $p$ as follows:

$$p \approx 1 - \frac{\int_1^{\text{IR} \cdot n} \frac{1}{x^\beta} dx}{\int_1^n \frac{1}{x^\beta} dx} \qquad (7)$$

Evaluating these integrals, we get:

$$p \approx 1 - \frac{\frac{1-(\text{IR} \cdot n)^{1-\beta}}{1-\beta}}{\frac{1-n^{1-\beta}}{1-\beta}} = 1 - \frac{1 - (\text{IR} \cdot n)^{1-\beta}}{1 - n^{1-\beta}} \qquad (8)$$

Substituting this expression for $p$ back into the complexity formula $C(\text{IR})$, we obtain:

$$C(\text{IR}) = \log(\text{IR} \cdot n) + \left(1 - \frac{1 - (\text{IR} \cdot n)^{1-\beta}}{1 - n^{1-\beta}}\right) \cdot \log n \qquad (9)$$

To find the optimal IR that minimizes the complexity $C(\text{IR})$, we take the derivative of $C(\text{IR})$ with respect to IR:

$$\frac{dC}{d\text{IR}} = \frac{1}{\text{IR}} + \frac{\log n \cdot (1 - \beta) \cdot n \cdot (\text{IR} \cdot n)^{-\beta}}{1 - n^{1-\beta}} \qquad (10)$$

Setting the derivative equal to zero to find the critical points:

$$\frac{1}{\text{IR}} + \frac{\log n \cdot (1 - \beta) \cdot n \cdot (\text{IR} \cdot n)^{-\beta}}{1 - n^{1-\beta}} = 0 \qquad (11)$$

After algebraic manipulation, the optimal IR can be expressed as:

$$\text{IR} = \left(\frac{n^{1-\beta} - 1}{(1 - \beta) \log n \cdot n^{1-\beta}}\right)^{\frac{1}{1-\beta}} \qquad (12)$$

For example, when $n = 1,000,000$ and set $\beta = 1.2$ as per Zipf's distribution of hot events in search engines [35], substituting these values into the formula gives a theoretical optimal IR of approximately 0.002. However, in practical scenarios, some adjustments are necessary. In real-world applications, it is often required to perform a certain amount of search within the full index to gather sufficient features and determine whether a node is of high frequency. This means that the full index's weight in the complexity calculation is relatively higher than that derived from the theoretical model. Consequently, based on practical experience and experimental observations, we have chosen an IR of 0.01 for our experiments, ensuring efficient and adaptive high-dimensional nearest neighbor search performance in dynamic query environments.

Solving this equation for IR yields the optimal index ratio that balances the trade-off between the hot index and full index search complexities. This analysis provides a theoretical foundation for dynamically adjusting the size of the hot index in response to varying query distributions and user preferences, ensuring efficient and adaptive high-dimensional nearest neighbor search performance.

## 5 EXPERIMENT

In this section, we evaluate the performance of DQF and conduct comparative evaluations with existing ANNS methods.

**Table 3: Dataset Information**

| Dataset | Dimension | Intrinsic Dimension | Dataset Size |
|---------|-----------|---------------------|--------------|
| SIFT1M | 128 | 12.9 | 1,000,000 |
| GIST1M | 960 | 29.1 | 1,000,000 |
| Glove | 100 | 20.9 | 1,183,514 |
| Crawl | 300 | 15.7 | 1,989,995 |

**Table 4: Evaluation Parameters**

| Parameter | Value |
|-----------|-------|
| Neighbor number $k$ | 1, 5, **10**, 20, 50 |
| Index Ratio $IR$ | 0.001, **0.005**, 0.01, 0.05, 0.1 |
| Stop Judgement Frequency $Freq$ | 20, **50**, 100, 200, 500 |
| Decision Tree Depth | 2, 5, **10**, 20, 50 |
| Add Step | **0**, 100, 200, 300, 400 |

## 5.1 Experimental Setup

*5.1.1 Experimental Datasets.* In this experiment, we selected four publicly available real-world datasets to comprehensively evaluate the performance of the proposed algorithm: SIFT1M[1], GIST1M[1], Glove[2], and Crawl[3]. The detailed dataset information is shown in Table 3. SIFT1M and GIST1M are datasets containing 1 million image feature vectors and are widely used in the field of image retrieval. The Glove dataset is a word vector dataset trained on 2 billion articles from Twitter text and is suitable for text semantic analysis. Crawl is a word vector collection constructed from text data crawled from the web and can represent the semantic features of large-scale web corpora. By conducting experiments on various data types such as images and text, we verified the generalization ability of our algorithm in different application scenarios.

*5.1.2 Experimental Settings.* The experiments were conducted on a computer with Intel(R) Xeon(R) Silver 4310 CPU@2.10GHz, 128GB of memory, and a 1TB hard drive. During the index construction phase, we used 24-core parallel processing to accelerate the construction. However, in the query phase, we switched to single-threaded operation to ensure accurate measurement of the search performance. All indexes were kept in memory for fast access.

For each dataset, we divided it into a training set and a test set in a ratio of 9 : 1. To simulate user query preferences, we generated query requests following the Zipf distribution with the distribution characteristics of hot events in China ($\beta = 1.2$) [35] as historical searches. Then, we generated 1,000 queries from the test set to calculate the recall rate. In the experiments, we set $K = 10$, which means each query retrieves 10 nearest neighbors as the final answer.

*5.1.3 Comparison Algorithms.* In this experiment, we compared our approach with the currently best-performing graph index algorithms, which are listed below:

**Table 5: Index Construction Time (the plus sign indicates the hot index time)**

| | HNSW | NSG | NSSG | DQF |
|---|------|-----|------|-----|
| SIFT1M | 116s | 137s | 79s | 79s+**1s** |
| GIST1M | 1222s | 1401s | 649s | 649s+**8s** |
| Crawl | 880s | 2487s | 970s | 970s+**12s** |
| Glove | 667s | 566s | 456s | 456s+**8s** |

**Table 6: Index Sizes (the plus sign indicates the hot index size)**

| | HNSW | NSG | NSSG | DQF |
|---|------|-----|------|-----|
| SIFT1M | 186MB | 104MB | 137MB | 137MB+**1MB** |
| GIST1M | 255MB | 77MB | 123MB | 123MB+**1MB** |
| Crawl | 302MB | 89MB | 161MB | 161MB+**2MB** |
| Glove | 220MB | 58MB | 90MB | 90MB+**1MB** |

- **HNSW** [27]: This algorithm constructs a multi-layer index based on a hierarchical navigable small-world graph, improving search efficiency. It starts from the top layer and searches downward layer by layer to accurately locate the nodes closest to the query point.
- **NSG** [14]: By utilizing a monotonic relative adjacency graph, NSG prunes non-monotonic edges, reducing memory usage. It efficiently finds target nodes without backtracking, making it effective for handling large-scale vector data.
- **NSSG** [13]: This algorithm introduces a satellite system graph pruning strategy, which ensures even distribution of outgoing edges of nodes, reduces the complexity of index construction, and enhances search performance.

In the experiments, the construction parameters for the comparison algorithms were set according to the configuration in the NSSG paper [13]. Meanwhile, the parameters for the hot index and the full index were also kept consistent with NSSG. Table 4 lists some parameters and their corresponding values used in the experiment, with defaults shown in bold.

## 5.2 Construction Performance

In this section, we evaluate the construction performance of our algorithm in terms of index construction time and index size.

As shown in Table 5, our DQF algorithm demonstrates significant efficiency in constructing the hot index. The index construction time for DQF includes two parts: building the full initial graph index using NSSG and the additional time for constructing the hot index. The hot index construction time is very short, allowing quick adaptation to new query patterns without rebuilding the entire index. In contrast, algorithms like HNSW, NSG, and NSSG require substantial time to rebuild the entire index when query distributions change, increasing system maintenance and response times. DQF's rapid adaptation capability gives it an edge in dynamic query scenarios, enabling timely responses to query preference changes and maintaining efficient search performance.
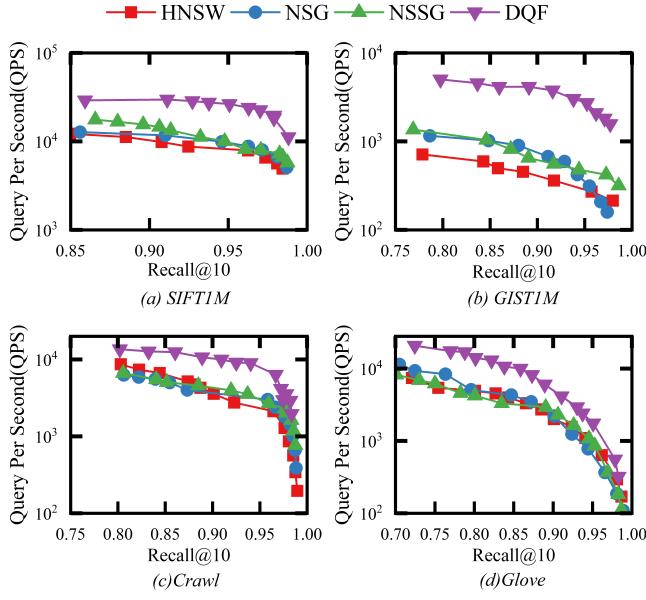
**Figure 5: Comparison of Search Performance**

Table 6 presents the index sizes of different algorithms across datasets. HNSW's multi-layer structure results in larger index sizes, while NSG's sparse graph construction yields the smallest index space. NSSG's index size lies between the two. Our DQF algorithm shares the same full graph index size as NSSG but adds a minimal hot index size (around 1MB). This makes DQF ideal for disk-memory hybrid storage scenarios [5, 21]: the full graph index can be stored on disk to save memory, while the hot index resides in memory for fast access. When query distributions change, only the in-memory hot index needs to be rebuilt, without altering the full graph index. This enhances system adaptability and query efficiency. In contrast, other algorithms typically require rebuilding the entire index when query preferences change, leading to higher storage and maintenance costs due to their larger index sizes.

### 5.3 Search Performance

*5.3.1 The Recall vs. Time.* Figure 5 illustrates the Recall@10 versus Query Per Second (QPS) performance of our Dual-Index Query Framework (DQF) compared to HNSW, NSG, and NSSG across four datasets. In the SIFT1M dataset, DQF demonstrates a 1.9× higher QPS than NSSG at a 90% recall rate and 1.5× higher at a 95% recall rate. For GIST1M, the QPS of DQF is 7.5× higher than NSSG at 90% recall and 5.7× higher at 95% recall. In the Crawl dataset, DQF outperforms NSSG by 2.4× at 90% recall and 2.8× at 95% recall. On the Glove dataset, DQF shows a 1.8× higher QPS than NSSG at 90% recall and 2.0× higher at 95% recall. These results clearly indicate that our DQF framework achieves significantly better query performance across different recall rates compared to NSSG, highlighting its efficiency and effectiveness in handling high-dimensional nearest neighbor search tasks.

*5.3.2 Effect of K.* Figure 6 presents the impact of varying $K$ values on both Query Per Second (QPS) and Recall across four datasets. In

terms of QPS, our Dual-Index Query Framework (DQF) generally outperforms HNSW, NSG, and NSSG in most scenarios. However, there are some exceptions. For instance, in the Glove dataset, the Recall is relatively low compared to other datasets. This may be attributed to the inherent complexity of the Glove dataset, which requires a larger $L$ value to achieve a high Recall. While our early termination strategy significantly enhances QPS, it may compromise precision slightly due to the dataset's complexity.

Regarding query performance, our framework consistently surpasses other methods in most cases. Notably, when $K = 50$, the search efficiency declines across all datasets. This is primarily because the hot index needs to accommodate more nodes as $K$ increases. Despite this, our proposed framework still outperforms the compared methods. It is worth mentioning that our framework performs much better with smaller $K$ values. This is because smaller $K$ values allow the hot index to more effectively focus on high-frequency queries, which aligns with the design philosophy of our dual-index structure. In summary, while there is room for improvement in handling larger $K$ values and complex datasets (e.g., Glove), our framework demonstrates superior overall performance and efficiency compared to other methods.

### 5.4 Parameter Analysis

*5.4.1 Effect of Index Ratio.* Figure 7 demonstrates the impact of different Index Ratios (IR) on the performance of the Dual-Index Query Framework across four datasets. From the results, it is evident that both excessively low and high IR values can negatively affect performance. When IR is too low (e.g., IR=0.001), the hot index may not encompass enough high-frequency nodes. This insufficiency forces the framework to frequently access the full index, even for high-frequency queries, thereby increasing the search time and reducing query throughput. Conversely, when IR is too high (e.g., IR=0.1), the hot index becomes overly large. This expansion dilutes the benefits of prioritizing high-frequency nodes and introduces unnecessary complexity to the hot index, which leads to decreased search efficiency.

A moderate IR maintains a compact hot index that sufficiently covers high-frequency nodes while avoiding the overhead of an excessively large index. This optimal IR setting allows the framework to efficiently handle the majority of high-frequency queries within the hot index, minimizing access to the full index and maximizing query performance.

*5.4.2 Impact of Depth and Frequency.* Figure 8 illustrates the impact of decision tree depth on the Dual-Index Query Framework's performance across four datasets. The results indicate that the framework's performance remains relatively stable across varying depths. When the depth is too low, the decision tree may lack the capacity to effectively distinguish between high-frequency and low-frequency queries, potentially reducing search efficiency. However, even at lower depths, the framework still performs adequately, likely due to the relatively simple decision boundaries needed for query differentiation.

As the depth increases, the decision tree gains more capacity to capture complex patterns in query features, refining routing decisions and improving search efficiency. Beyond a certain point, though, the marginal gains in performance diminish. This suggests
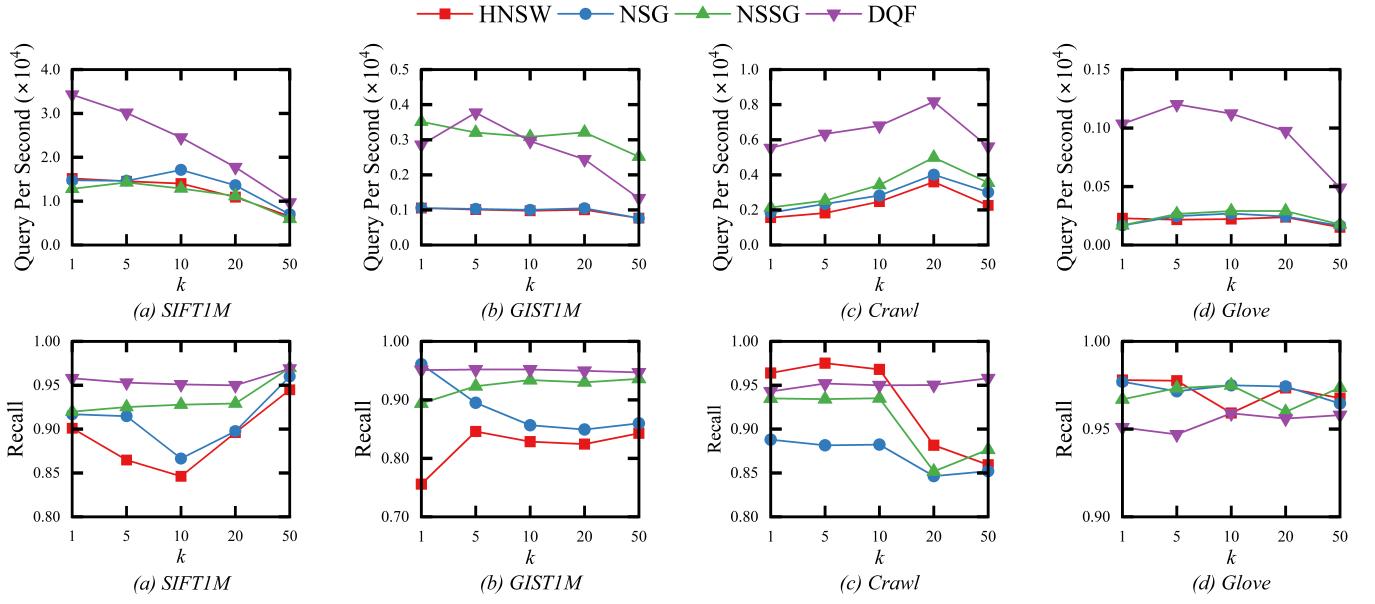
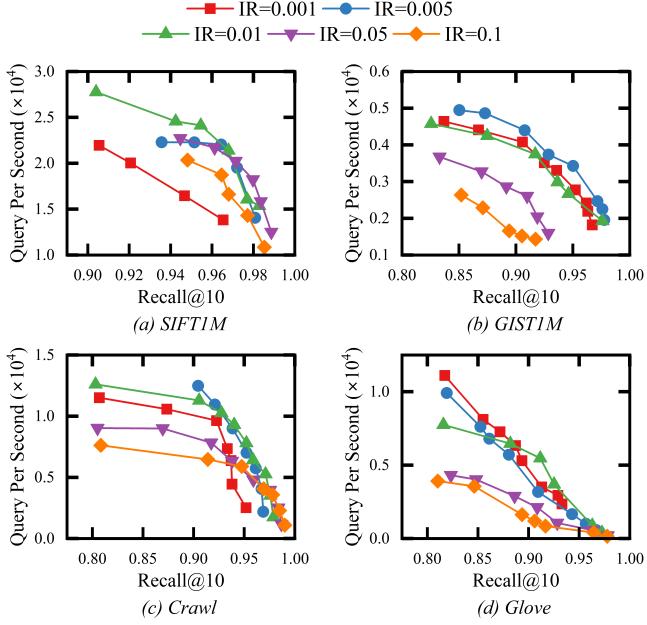Figure 6: Impact of $k$ on DQF

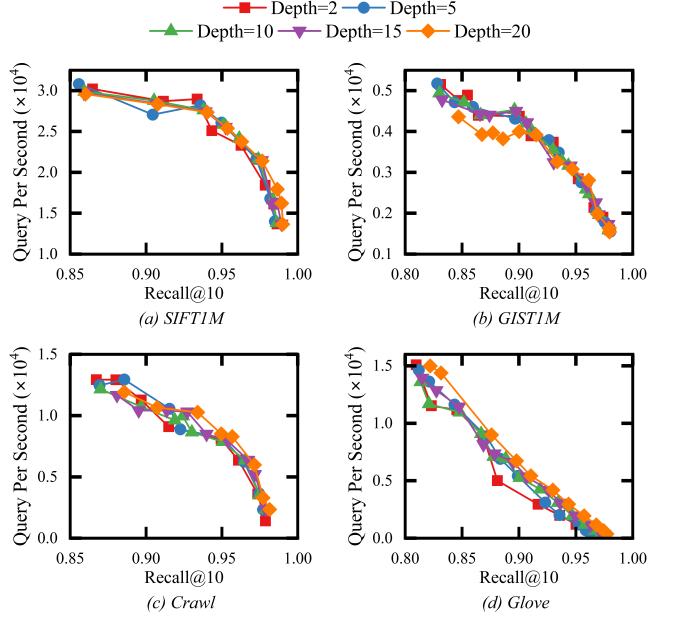

Figure 7: Effect of Index Ratio on DQF



Figure 8: Effect of Decision Tree Depth on DQF

that a moderate depth is sufficient for the decision tree to achieve effective routing without unnecessary complexity. The results show that setting the decision tree depth to a moderate level ensures efficient and effective search operations across different datasets.

Figure 9 illustrates the impact of varying decision tree evaluation gaps (*EvalGap*) on the Dual-Index Query Framework's performance across four datasets. *EvalGap* indicates the number of distance computations performed between successive decision tree

evaluations. The findings reveal that lowering *EvalGap* generally enhances performance. A smaller *EvalGap*, such as 20, results in more frequent decision tree invocations. This allows for timelier decisions to terminate low-frequency queries early, minimizing unnecessary computations and boosting query throughput. Conversely, a larger *EvalGap* reduces the frequency of decision tree calls, potentially leading to delayed termination of low-frequency queries and increased computational overhead.
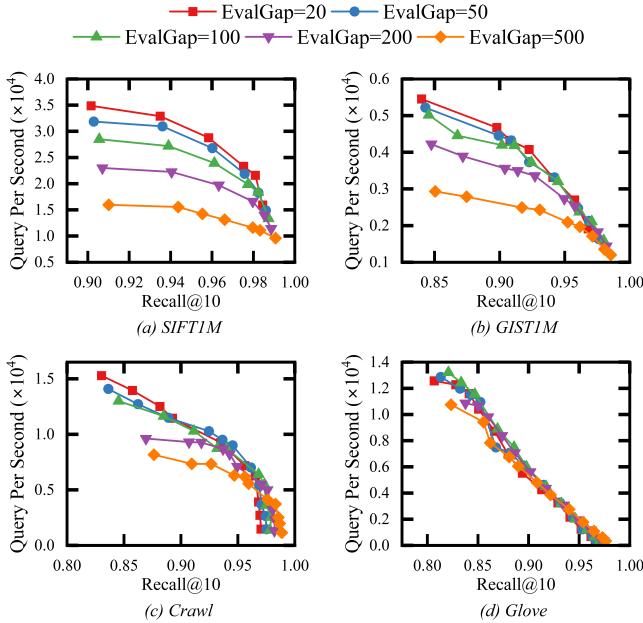
Figure 9: Effect of EvalGap on DQF



Figure 10: Add Step After Decision Tree Terminate

Counterintuitively, too small *EvalGap* can significantly boost the decision tree's computational overhead, overshadowing the time saved in distance calculations. However, the decision tree's computational cost is relatively small compared to the efficiency gains it provides in optimizing the search strategy. The results show that increasing the judgment frequency improves the framework's dynamic adaptation to query characteristics, leading to better resource allocation and search performance. Despite minor performance variations with different *EvalGap* settings, our proposed framework consistently outperforms the compared methods, underscoring the robustness and effectiveness of the decision tree-based dynamic search strategy.

*5.4.3 Add Step.* Figure 10 shows the effect of the *AddStep* parameter on the performance of the Dual-Index Query Framework across four datasets. The *AddStep* parameter determines the number of additional search steps executed after the decision tree triggers a termination. The results indicate that increasing the *AddStep* value generally has minimal impact on both query throughput and recall.

When the parameter *AddStep* is set to 0, the framework relies solely on the decision tree to determine when to terminate the search. The results show that even with *AddStep* = 0, the framework achieves strong performance, suggesting that the decision tree effectively identifies optimal termination points for most queries. Adding more steps (e.g., *AddStep* = 100, 200, etc.) does not significantly improve recall or query per second (QPS), indicating that the decision tree already makes accurate termination decisions without requiring additional search steps. This efficiency is likely due to the decision tree's ability to leverage query-specific features to refine the search strategy adaptively.

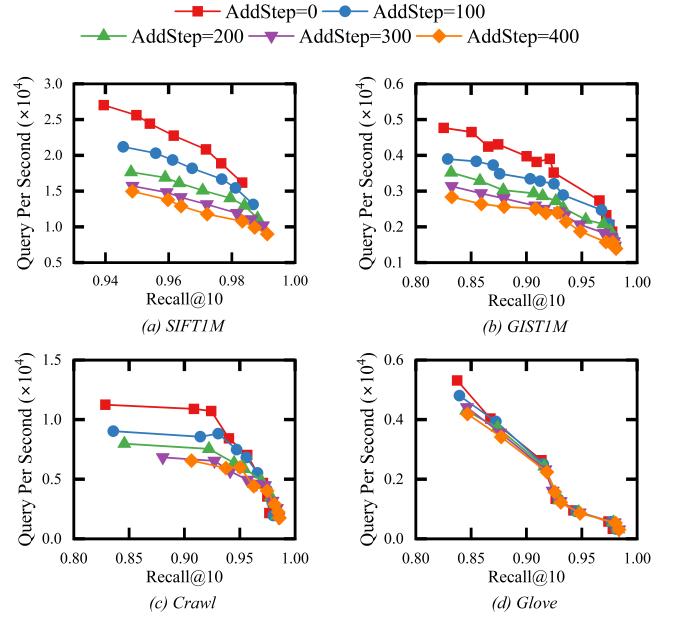The minimal gains from increasing the *AddStep* parameter imply that the decision tree's termination decisions are robust and that additional steps do not substantially enhance the search outcomes. This supports the effectiveness of the decision tree in dynamically adjusting the search process, reducing unnecessary computations while maintaining high search efficiency and recall.

## 6 CONCLUSION

In conclusion, we have presented a novel Dual-Index Query Framework that effectively addresses the challenges of high-dimensional nearest neighbor search by integrating user query preferences and temporal dynamics. Through a dual-layer index structure and a dynamic decision-tree-powered search strategy, our framework achieves significant performance improvements, offering faster query processing and high recall rates while adapting efficiently to dynamic query patterns. Experimental results demonstrate the effectiveness and practicality of our proposed framework, highlighting its potential for real-world applications where query preferences and timeliness are critical.

## REFERENCES

[1] Lada A Adamic and Bernardo A Huberman. 2000. Power-law distribution of the world wide web. *science* 287, 5461 (2000), 2115–2115.

[2] Lada A Adamic and Bernardo A Huberman. 2002. Zipf's law and the Internet. *Glottometrics* 3, 1 (2002), 143–150.

[3] Norbert Beckmann, Hans-Peter Kriegel, Ralf Schneider, and Bernhard Seeger. 1990. The R*-tree: An efficient and robust access method for points and rectangles. In *Proceedings of the 1990 ACM SIGMOD international conference on Management of data*. 322–331.

[4] Jon Louis Bentley. 1975. Multidimensional binary search trees used for associative searching. *Commun. ACM* 18, 9 (1975), 509–517.

[5] Qi Chen, Bing Zhao, Haidong Wang, Mingqin Li, Chuanjie Liu, Zengzhong Li, Mao Yang, and Jingdong Wang. 2021. Spann: Highly-efficient billion-scale approximate nearest neighborhood search. *Advances in Neural Information Processing Systems* 34 (2021), 5199–5212.

[6] Rihan Chen, Bin Liu, Han Zhu, Yaoxuan Wang, Qi Li, Buting Ma, Qingbo Hua, Jun Jiang, Yunlong Xu, Hongbo Deng, et al. 2022. Approximate nearest neighbor search under neural similarity metric for large-scale recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge*

*Management*. 3013–3022.

[7] Tingyang Chen, Cong Fu, Kun Wang, Xiangyu Ke, Yunjun Gao, Wenchao Zhou, Yabo Ni, and Anxiang Zeng. 2025. Maximum inner product is query-scaled nearest neighbor. *arXiv preprint arXiv:2503.06882* (2025).

[8] Aaron Clauset, Cosma Rohilla Shalizi, and Mark EJ Newman. 2009. Power-law distributions in empirical data. *SIAM review* 51, 4 (2009), 661–703.

[9] Arjen P de Vries, Nikos Mamoulis, Niels Nes, and Martin Kersten. 2002. Efficient k-NN search on vertically decomposed data. In *Proceedings of the 2002 ACM SIGMOD international conference on Management of data*. 322–333.

[10] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitansky, Robert Osazuwa Ness, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130* (2024).

[11] Ada Wai-chee Fu, Polly Mei-shuen Chan, Yin-Ling Cheung, and Yiu Sang Moon. 2000. Dynamic vp-tree indexing for n-nearest neighbor search given pair-wise distances. *The VLDB Journal* 9 (2000), 154–173.

[12] Cong Fu and Deng Cai. 2016. Efanna: An extremely fast approximate nearest neighbor search algorithm based on knn graph. *arXiv preprint arXiv:1609.07228* (2016).

[13] Cong Fu, Changxu Wang, and Deng Cai. 2021. High dimensional similarity search with satellite system graph: Efficiency, scalability, and unindexed query compatibility. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 8 (2021), 4139–4150.

[14] Cong Fu, Chao Xiang, Changxu Wang, and Deng Cai. 2019. Fast approximate nearest neighbor search with the navigating spreading-out graph. *Proceedings of the VLDB Endowment* 12, 5 (2019), 461–474.

[15] Tiezheng Ge, Kaiming He, Qifa Ke, and Jian Sun. 2013. Optimized product quantization for approximate nearest neighbor search. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2946–2953.

[16] Phillipa Gill, Martin Arlitt, Zongpeng Li, and Anirban Mahanti. 2007. Youtube traffic characterization: a view from the edge. In *Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*. 15–28.

[17] Aristides Gionis, Piotr Indyk, Rajeev Motwani, et al. 1999. Similarity search in high dimensions via hashing. In *Vldb*, Vol. 99. 518–529.

[18] Long Gong, Huayi Wang, Mitsunori Ogihara, and Jun Xu. 2020. iDEC: indexable distance estimating codes for approximate nearest neighbor search. *Proceedings of the VLDB Endowment* 13, 9 (2020).

[19] Qiang Huang, Jianlin Feng, Yikai Zhang, Qiong Fang, and Wilfred Ng. 2015. Query-aware locality-sensitive hashing for approximate nearest neighbor search. *Proceedings of the VLDB Endowment* 9, 1 (2015), 1–12.

[20] Piotr Indyk and Rajeev Motwani. 1998. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*. 604–613.

[21] Suhas Jayaram Subramanya, Fnu Devvrit, Harsha Vardhan Simhadri, Ravishankar Krishnaswamy, and Rohan Kadekodi. 2019. Diskann: Fast accurate billion-point nearest neighbor search on a single node. *Advances in neural information processing Systems* 32 (2019).

[22] Herve Jegou, Matthijs Douze, and Cordelia Schmid. 2010. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence* 33, 1 (2010), 117–128.

[23] Jie Li, Haifeng Liu, Chuanghua Gui, Jianyu Chen, Zhenyuan Ni, Ning Wang, and Yuan Chen. 2018. The design and implementation of a real time visual search system on JD E-commerce platform. In *Proceedings of the 19th International Middleware Conference Industry*. 9–16.

[24] Fabrizio Lillo and Salvatore Ruggieri. 2019. Estimating the total volume of queries to Google. In *The World Wide Web Conference*. 1051–1060.

[25] Kejing Lu, Hongya Wang, Wei Wang, and Mineichi Kudo. 2020. VHP: approximate nearest neighbor search via virtual hypersphere partitioning. *Proceedings of the VLDB Endowment* 13, 9 (2020), 1443–1455.

[26] Yury Malkov, Alexander Ponomarenko, Andrey Logvinov, and Vladimir Krylov. 2014. Approximate nearest neighbor algorithm based on navigable small world graphs. *Information Systems* 45 (2014), 61–68.

[27] Yu A Malkov and Dmitry A Yashunin. 2018. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence* 42, 4 (2018), 824–836.

[28] John Paparrizos, Ikraduya Edian, Chunwei Liu, Aaron J Elmore, and Michael J Franklin. 2022. Fast adaptive similarity search through variance-aware quantization. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*. IEEE, 2969–2983.

[29] J Ben Schafer, Dan Frankowski, Jon Herlocker, and Shilad Sen. 2007. Collaborative filtering recommender systems. In *The adaptive web: methods and strategies of web personalization*. Springer, 291–324.

[30] Ján Suchal and Pavol Návrat. 2010. Full text search engine as scalable k-nearest neighbor recommendation system. In *Artificial Intelligence in Theory and Practice III: Third IFIP TC 12 International Conference on Artificial Intelligence, IFIP AI 2010, Held as Part of WCC 2010, Brisbane, Australia, September 20-23, 2010. Proceedings 3*. Springer, 165–173.

[31] Godfried T Toussaint. 1980. The relative neighbourhood graph of a finite planar set. *Pattern recognition* 12, 4 (1980), 261–268.

[32] Mengzhao Wang, Xiaoliang Xu, Qiang Yue, and Yuxiang Wang. 2021. A comprehensive survey and experimental comparison of graph-based approximate nearest neighbor search. *Proceedings of the VLDB Endowment* 14, 11 (2021), 1964–1978.

[33] Runhui Wang and Dong Deng. 2020. DeltaPQ: lossless product quantization code compression for high dimensional similarity search. *Proceedings of the VLDB Endowment* 13, 13 (2020), 3603–3616.

[34] Roger Weber, Hans-Jörg Schek, and Stephen Blott. 1998. A quantitative analysis and performance study for similarity-search methods in high-dimensional spaces. In *VLDB*, Vol. 98. 194–205.

[35] Yingfan Xu and Mingliang Shi. 2018. Research on zipf's law of hot events in search engines. In *WHICEB 2018 Proceedings*, Vol. 48. 189–195.

[36] Xizhe Yin, Chao Gao, Zhijia Zhao, and Rajiv Gupta. 2025. PANNS: Enhancing graph-based approximate nearest neighbor search through recency-aware construction and parameterized search. In *Proceedings of the 30th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming*. 369–381.

[37] Ziqi Yin, Jianyang Gao, Pasquale Balsebre, Gao Cong, and Cheng Long. 2025. DEG: Efficient hybrid vector search using the dynamic edge navigation graph. *Proceedings of the ACM on Management of Data* 3, 1 (2025), 1–28.

[38] Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng, Fangcheng Fu, Ling Yang, Wentao Zhang, Jie Jiang, and Bin Cui. 2024. Retrieval-augmented generation for ai-generated content: A survey. *arXiv preprint arXiv:2402.19473* (2024).

[39] Wayne Xin Zhao, Jing Liu, Ruiyang Ren, and Ji-Rong Wen. 2024. Dense text retrieval based on pretrained language models: A survey. *ACM Transactions on Information Systems* 42, 4 (2024), 1–60.