# Distribution-Aware Exploration for Adaptive HNSW Search

Chao Zhang
University of Waterloo
Waterloo, Canada
chao.zhang@uwaterloo.ca

Renée J. Miller
University of Waterloo
Waterloo, Canada
rjmiller@uwaterloo.ca

## Abstract

Hierarchical Navigable Small World (HNSW) is widely adopted for approximate nearest neighbor search (ANNS) for its ability to deliver high recall with low latency on large-scale, high-dimensional embeddings. The exploration factor, commonly referred to as ef, is a key parameter in HNSW-based vector search that balances accuracy and efficiency. However, existing systems typically rely on manually and statically configured ef values that are uniformly applied across all queries. This results in a distribution-agnostic configuration that fails to account for the non-uniform and skewed nature of real-world embedding data and query workloads. As a consequence, HNSW-based systems suffer from two key practical issues: (i) the absence of recall guarantees, and (ii) inefficient ANNS performance due to over- or under-searching. In this paper, we propose Adaptive-ef (Ada-ef), a data-driven, update-friendly, query-adaptive approach that dynamically configures ef for each query at runtime to approximately meet a declarative target recall with minimal computation. The core of our approach is a theoretically grounded statistical model that captures the similarity distribution between each query and the database vectors. Based on this foundation, we design a query scoring mechanism that distinguishes between queries requiring only small ef and those that need larger ef to meet a target recall, and accordingly assigns an appropriate ef to each query. Experimental results on real-world embeddings produced by state-of-the-art Transformer models from OpenAI and Cohere show that, compared with state-of-the-art learning-based adaptive approaches, our method achieves the target recall while avoiding both over- and under-searching, reducing online query latency by up to 4×, offline computation time by 50×, and offline memory usage by 100×.

## 1 Introduction

The rapid advancement of embedding models [67], particularly Transformer-based architectures [32, 34, 79, 87], has revolutionized the representation and processing of data. These models encode diverse types of information, such as text [56] and images [77], into high-dimensional semantic embeddings that capture contextual meaning in vector space. To compare these embeddings effectively, inner product and cosine similarity [69, 77, 79, 99] have become standard, due to their ability to reflect semantic alignment between modern semantic embeddings. Building on this foundation, vector indexing techniques are employed to efficiently retrieve the most relevant embeddings from large collections, making them critical components in real-time semantic search and retrieval systems.

Various types of vector indexes have been developed, including tree-based [18, 25, 60, 63, 81], hashing-based [45, 49, 50, 101], quantization-based [17, 55, 73], and graph-based [21, 21, 30, 33, 40–42, 47, 51, 53, 54, 61, 64, 65, 86] methods. Among these, graph-based indexes have emerged as the most effective for high-dimensional similarity search, leveraging proximity graphs to efficiently structure vector data. One of the most widely adopted graph-based solutions for vector search is Hierarchical Navigable Small World (HNSW) [65]. It serves as a core indexing technique in vector search libraries like Faiss [35], and powers search engines such as Elasticsearch [84] and Lucene [10]. Additionally, it is the backbone of specialized vector databases, including Milvus [88], Pinecone [5], and Weaviate [6], and has also been integrated into relational databases such as PostgreSQL [9, 97], SingleStore [28], and DuckDB [43, 76], and graph databases like TigerGraph [62] and Kùzu [38].
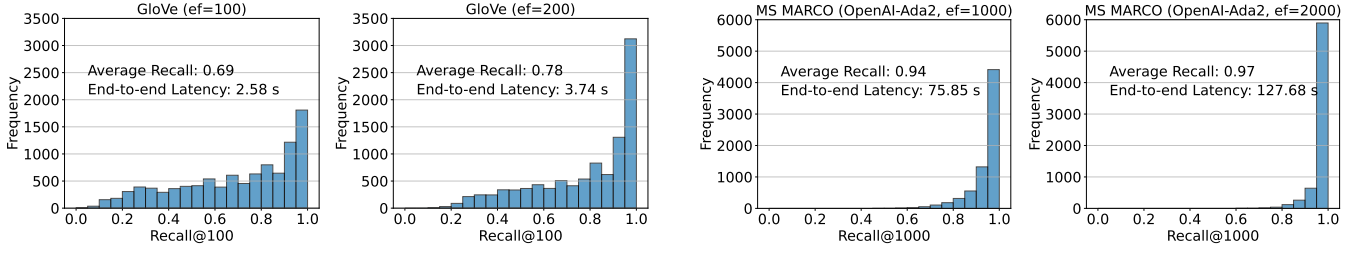
Given a set of data vectors $\mathbf{V}$, HNSW constructs a hierarchical proximity graph to serve as a vector index for approximate nearest neighbor search (ANNS) over $\mathbf{V}$. The number of nodes decreases logarithmically from the bottom layer to the top level, with the bottom level containing all $|\mathbf{V}|$ nodes, each representing a data vector. Edges in the HNSW graph encode the similarity relationship between vectors. During query time, the search traverses the hierarchical graph from top to bottom to identify an entry point at the base layer. Starting from this point, the search proceeds using a best-first strategy guided by a fixed-size priority queue, whose size is a user-defined parameter, commonly referred to as ef or efSearch, short for *exploration factor*. The search terminates when the furthest element currently in the queue is closer to the query than the closest candidate node to be visited, indicating that no better candidate is likely to be found.

The parameter ef controls the trade-off between query latency and result accuracy. A larger ef typically yields a higher recall at the cost of increased query time, whereas a smaller ef results in faster queries but may compromise accuracy. The default setting of ef in existing systems [36] is either a static ef value, *e.g.*, ef=40 in PGVector [9] and ef=16 in Faiss [35], or setting ef to $k$ for Top-$k$ ANNS, where the goal is to retrieve the $k$ nearest vectors to a given query, *e.g.*, OpenSearch (Lucene) [8].

Once configured, the ef value remains fixed and is applied uniformly across the entire query workload, under the expectation that it will consistently provide high recall. However, this is not guaranteed so to compensate, users typically set ef to a sufficiently large value, often based on empirical heuristics or subjective estimation.

**Accepted for publication in SIGMOD 2026**

**Figure 1: Recall distribution of HNSW search. GloVe [75]: 1.8M vectors (100D), 10K queries, Top-100 ANNS with `ef` = 100 and 200. MS MARCO (OpenAI Ada-002 embeddings) [11]: 8.8M vectors (1536D), 6.9K queries, Top-1000 ANNS with `ef` = 1000 and 2000.**

However, this static `ef` configuration introduces two key issues in production environments: (i) the absence of recall guarantees; and (ii) inappropriate `ef` values for individual queries, leading to either over-searching or under-searching.

More fundamentally, applying a fixed `ef` uniformly across queries is structurally misaligned with the inherently non-uniform and skewed nature of semantic embedding spaces. Semantic embeddings, including word, sentence, and image embeddings, are widely known to exhibit highly skewed, non-uniform distributions. In word embeddings, this skew is often tied to frequency bias: frequent words dominate the embedding space with larger norms and concentrated directional variance, a phenomenon shown to degrade similarity search due to the emergence of hub-points that appear as nearest neighbors to many unrelated vectors [68, 98]. Similar skew is observed in sentence embeddings from transformer-based models like BERT and RoBERTa, which tend to produce anisotropic vector spaces where embeddings cluster along dominant directions, limiting their discriminative power [37, 58, 70]. Cross-modal embeddings, such as CLIP's joint image–text vectors, also suffer from skew and hubness: normalized vectors lie in narrow cones, with retrieval dominated by a small set of hub images or texts [26, 85].

**Motivating example**. To illustrate these issues, we built HNSW indexes using `HNSWlib` [2] for two datasets (see Table 1 for dataset statistics) and conduct Top-$k$ ANNS with static `ef` values. The results in Figure 1 highlight several important observations. First, different datasets require different `ef` to achieve high average recall. For example, setting `ef` to $k$ yields relatively high recall in MS MARCO (average recall is 94%), but performs less effectively in GloVe (average recall is only 69%). Second, a significant portion of the queries exhibit recall substantially higher or lower than the average, indicating that average recall alone does not fully capture actual search performance. While increasing `ef` from $k$ to $2k$ can improve average recall, it also leads to over-searching for queries that already achieve high recall. In the GloVe dataset, approximately 1800 queries attain recall close to 1.0 with `ef`=$k$; increasing `ef` to $2k$ offers little improvement for these queries but adds unnecessary computation overhead. A similar trend is observed in the MS MARCO dataset, where over 4000 queries already perform well with `ef`= $k$, making further increases redundant. Additionally, even though the average recall in MS MARCO reaches 0.97 with `ef`= $2k$, approximately 1000 queries still have recall below 0.9. In the GloVe dataset, a significant number of queries fall well below the average.
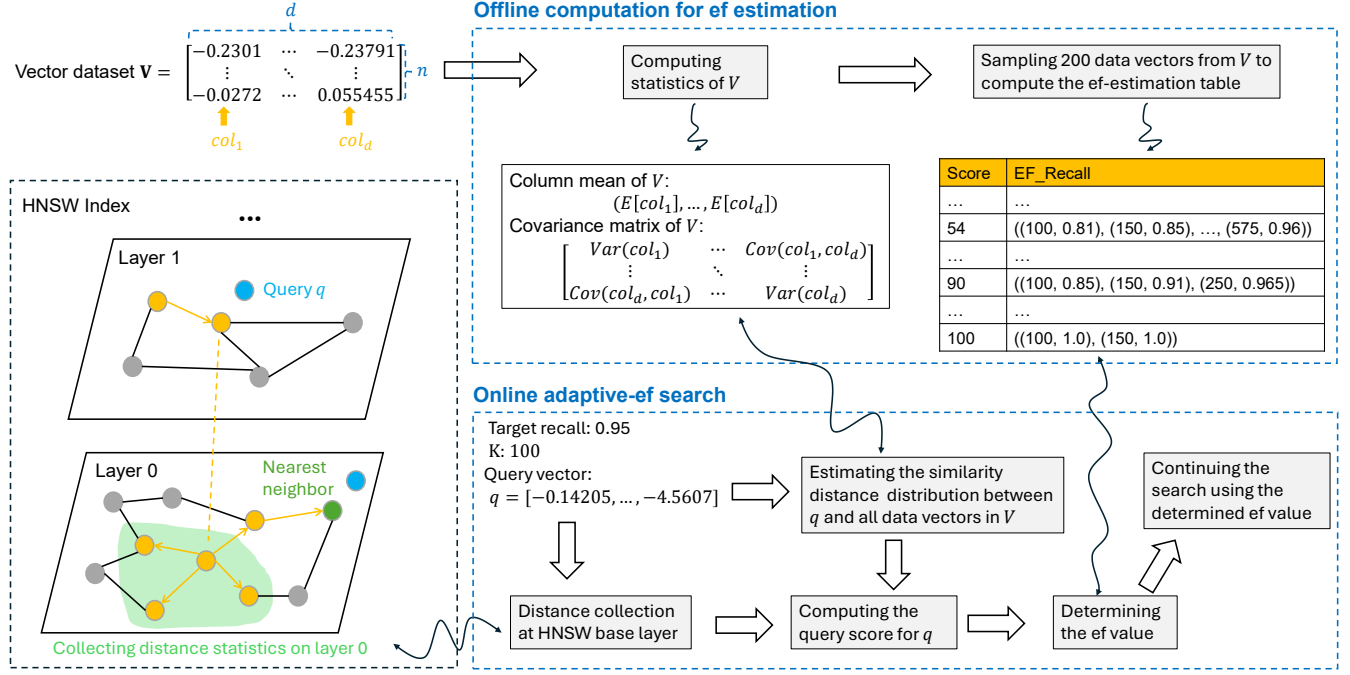
In real-world production environments, selecting an appropriate `ef` value to meet a target recall is highly challenging. Despite its practical importance, the problem of adaptively configuring `ef` to meet recall requirements has received little attention. Most benchmarking efforts evaluate indexing methods by sweeping across various `ef` values to generate latency-recall trade-off curves for the fixed query load of the benchmarks. While useful for comparing the relative performance of indexing structures, these evaluations do not offer a practical solution for configuring `ef` in real-world deployments. Specifically, it is infeasible to iteratively increase `ef` to reach a desired recall level for online query processing, as doing so would require repeated vector search and computing the ground truth for incoming query vectors. Moreover, although setting a large `ef` can improve overall recall, it also incurs substantial computation overhead and may still fail to achieve high recall across all queries, as demonstrated in the experiments.

In this paper, we investigate the following two key questions:
(i) *Can we configure* `ef` *to meet a user-specified target recall?*
(ii) *Can we assign adaptive* `ef` *values on a per-query basis, allocating a larger* `ef` *to queries requiring more computation to achieve the target recall, and a smaller* `ef` *to those for which a lower value suffices?*

To address the above challenges, we propose **Adaptive-ef (`Ada-ef`)**, a query-adaptive approach for `ef` configuration, *i.e.*, configuring an `ef` value for each individual query vector to approximately achieve a target recall at runtime. `Ada-ef` is lightweight, update-friendly, and purely rule-based, which allows it to generalize effectively across different datasets. Figure 2 presents an overview of our approach. The core idea is as follows: Given a query vector, we first perform the standard HNSW search to reach the base layer and obtain the entry point. The subsequent online search process consists of two phases: (i) Distance collection; and (ii) Search with a deterministic `ef` value. In phase (i), the algorithm explores the graph and records the distances between the query vector and the visited nodes. Based on these distances collected, we design an `ef`-estimator that takes as input the query vector, the collected distance list, and a user-defined target recall, and outputs an `ef` value intended to achieve the target recall for the given query, which is used in phase (ii).

The `ef`-estimator is the core component of our adaptive approach. Its design is based on our theoretical foundation on the similarity distance distribution between a query vector $q$ and a vector database: *Cosine similarity and inner product among high-dimensional learned embeddings often approximate a Gaussian (normal) distribution.* More importantly, the distribution can be estimated efficiently

**Figure 2: Overview of `Ada-ef`. It consists of two stages: offline and online computations. In the offline stage, dataset-level statistics are computed (§ 5), followed by the construction of an `ef`-estimation table (§ 6). In the online stage, the search follows the standard HNSW process until the base layer, where adaptive-`ef` search begins (§ 4): (1) Distance collection: exploring a limited number of nodes to compute a query score using offline statistics; (2) Search with estimated `ef`: using the score to select an `ef` from the estimation table to approximately reach a declarative target recall.**

at runtime. Combined with the distances collected during phase (i) of the search process, the estimated distribution allows us to distinguish between queries that can achieve high recall with a small `ef` and those that require a large `ef` to achieve a target recall.

Empirically, we demonstrate that our `ef`-adaptive approach is able to achieve the user-defined target recall across query workloads. Additionally, it can improve the recall of queries that suffer from under-searching under a static `ef` configuration, *i.e.*, the tail cases shown in Figure 1. Importantly, although our method involves collecting distance statistics and estimating `ef` for each query on the fly, the overall end-to-end query latency can be lower than that of the original HNSW search. This is attributed to the adaptive nature of our approach, which allocates computational effort more efficiently across queries, thereby mitigating the issue of over-searching. We further show that, while our method requires an offline preprocessing step, the associated computational cost and storage overhead are negligible compared to the indexing overhead of HNSW. Last but not least, by simulating a real-world production scenario, we demonstrate that our approach can effectively adapt to dynamic workloads involving both insertion and deletion, and that the corresponding offline updates are lightweight.

We make the following two major contributions in this paper:

- We introduce `Ada-ef`, a distribution-aware approach that adaptively sets `ef` for each query at runtime to approximately meet a user-specified target recall, while improving

efficiency by avoiding both under- and over-searching. To the best of our knowledge, this is the first work that addresses `ef` configuration for HNSW search.

- We establish a theoretical foundation for the distribution of similarity distances between a query vector and a vector database. We demonstrate the practical applicability of this foundation in the context of searching high-dimensional learned embeddings.

## 2 Related Work

We position our work within existing work. For a more comprehensive discussion, we refer readers to prior surveys [22, 52, 72, 91, 92].

Various indexing techniques have been proposed (see §1 for detail). Among them, HNSW [65] stands out as the most influential and widely adopted approach. HNSW has been extensively adopted (see §1 for detail). Although graph-based indexes [21, 21, 30, 33, 40–42, 47, 51, 53, 54, 61, 64, 65, 86] differ in their indexing algorithms for graph construction, their online search algorithms remain fundamentally the same across all methods [22, 92], such as NSG [42], Vamana [53], and SPANN [31]. A few studies have explored ways to further optimize the search process. These include leveraging intra-query parallelism to accelerate search execution [74] and approximating distance computations to improve efficiency while maintaining accuracy [29]. Several recent studies have also explored the

impact of hardware architectures on HNSW performance, including GPU acceleration [71], CPU-based optimizations [89], hybrid CPU-GPU approaches [83], and SSD-based indexing [90].

The key distinction between our work and all prior graph-based approaches lies in the use of a fixed versus an adaptive exploration factor in the best-first search algorithm, *i.e.*, the ef parameter. Existing methods rely on a fixed ef value, manually pre-configured by the user and applied uniformly across all queries during search. This static configuration lacks recall approximation and often leads to inefficiencies, *i.e.*, over-searching for some queries and under-searching for others, as shown in our motivation examples and experiments. In contrast, we introduce Ada-ef, the first distribution-aware approach for adaptively configuring ef on a per-query basis at runtime. By leveraging dataset- and query-specific statistics, this adaptive strategy enables a more fine-grained trade-off between accuracy and efficiency, and is particularly well-suited to real-world workloads, where data and query distributions are often non-uniform and highly skewed. We focus on HNSW, given its widespread adoption in practice. Our method operates on existing pre-built indexes without requiring any index modifications, enabling seamless integration into deployed infrastructures. Furthermore, the proposed adaptive strategy is orthogonal to both algorithmic and hardware-level optimizations, making it compatible with and complementary to existing acceleration techniques.

Early termination in ANNS has been explored through both learning-based and heuristic approaches. Learned Adaptive Early Termination (LAET) [59] predicts per-query stopping conditions using Gradient Boosting Decision Trees [39] and runtime features. Tao [95] is a static-feature-only framework that replaces runtime features with the Local Intrinsic Dimension (LID) [16]. DARTH [27] extends LAET by integrating a recall predictor within HNSW to achieve declarative recall via prediction intervals. In contrast, Patience in Proximity (PiP) [82] is a simple saturation-based heuristic that halts HNSW search once candidate improvements plateau. Unlike PiP, which requires defining a fixed ef value, LAET, Tao, and DARTH automatically determine termination points without preset ef parameters, relying instead on model predictions. Among these methods, DARTH represents the state-of-the-art and is the only one supporting declarative recall, offering the interface as Ada-ef.

Compared to LAET, Tao, DARTH, and PiP, Ada-ef is entirely rule-based and has a novel design based on our foundation on similarity distribution (Theorem 5.2). It estimates the ef value required to approximately reach a target recall. Consequently, Ada-ef incurs minimal offline computation (§7.3) and supports efficient incremental updates of insertion and deletion (§6.3). More importantly, Ada-ef achieves robust query performance, reaching target recall across datasets, reducing latency, and improving recall on the hardest queries, as demonstrated in our experimental evaluation.

VBASE [100] incorporates an early termination mechanism for hybrid search combining vector similarity and attribute filtering, enabling simultaneous evaluation of both components. However, its benefit is limited in pure vector search, where it still depends on a fixed exploration parameter $E$ (equivalent to ef in HNSW). Extending our query-adaptive ef approach to hybrid search remains a promising direction for future work.

## 3 Preliminary

A vector dataset $\mathbf{V} = (\mathbf{v_1}, ..., \mathbf{v_n})$ is a collection of $n$ vectors, where each vector $\mathbf{v_i} = (v_1, ..., v_d)$ has dimensionality $d$. Each $\mathbf{v_i}$ is referred to as a data vector, which typically represents an embedded form of raw data, such as an image or a text passage. We represent the entire vector dataset $\mathbf{V}$ as an $n \times d$ matrix. The similarity between a query vector $\mathbf{q}$ and a data vector $\mathbf{v}$ is measured using a distance function $dist(\mathbf{q}, \mathbf{v})$, where a smaller value indicates greater similarity. For instance, cosine distance is a common choice for $dist$ in the case of learned vector embeddings in modern AI systems.

*Definition 3.1 (Full Distance List).* Given a query vector $\mathbf{q}$ and a dataset $\mathbf{V} = (\mathbf{v_1}, ..., \mathbf{v_n})$, the Full Distance List (FDL) of $\mathbf{q}$ with respect to $\mathbf{V}$ is $FDL(\mathbf{q}, \mathbf{V}) = (dist(\mathbf{q}, \mathbf{v_1}), ..., dist(\mathbf{q}, \mathbf{v_n}))$.

Given a query vector $\mathbf{q}$ and a vector dataset $\mathbf{V}$, the objective of nearest neighbor search (NNS) is to find a set $\mathcal{R}$ of $k$ vectors in $\mathbf{V}$ that are closest to $\mathbf{q}$ under a given distance function $dist$. Formally, $|\mathcal{R}| = k$, and for any vector $\mathbf{v} \in \mathcal{R}$, there does not exist a vector $\mathbf{v'} \in \mathbf{V} \setminus \mathcal{R}$ such that $dist(\mathbf{q}, \mathbf{v}) > dist(\mathbf{q}, \mathbf{v'})$.

When both $n$ (the number of vectors) and $d$ (the dimensionality) are large, as is common in modern applications (*e.g.*, millions of vectors with 768 dimensions generated by BERT-based embedding models), computing exact NNS (*i.e.*, computing FDL) becomes prohibitively expensive. To address this scalability issue, practical systems typically employ ANNS, which computes an approximate result set $\tilde{\mathcal{R}}$ that aims to closely match the true nearest neighbor set $\mathcal{R}$. The accuracy of ANNS is commonly evaluated using the recall metric, defined as $Recall@k = \frac{|\mathcal{R} \cap \tilde{\mathcal{R}}|}{k}$.

Vector indexes are specialized data structures designed to accelerate ANNS by scanning only a subset of data vectors in $\mathbf{V}$ to compute an approximate result set $\tilde{\mathcal{R}}$. In this paper, we focus on ANNS using HNSW, the most widely adopted vector indexes in practice. Specifically, we consider scenarios where an HNSW index has already been constructed and focus on searching the index.

The key parameter in HNSW search is the size of the priority queue that maintains the search results at the base layer, commonly denoted as ef. This parameter controls the trade-off between recall and search latency, as illustrated in Figure 1. In practice, ef is typically configured based on subjective estimation or through iterative empirical tuning. Once selected, the same ef value is applied uniformly across the entire query workload, resulting in the absence of recall guarantees and inefficient search due to under- and over-searching. To overcome these issues, we study the problem of query-adaptive ef configuration defined below.

*Definition 3.2.* Given a HNSW index, a query vector $\mathbf{q}$, and a user-specified expected recall, the adaptive ef for $\mathbf{q}$ is defined as the minimum ef value such that the search result of $\mathbf{q}$ over the HNSW index achieves the expected recall.

We propose a distribution-aware, data-driven approach to estimate adaptive ef per query, which can be seamlessly integrated into an existing HNSW index.

## 4 Adaptive EF Search

We introduce Ada-ef, an alternative to standard HNSW search at the base-layer proximity graph. The key distinction is that Ada-ef

does not require a user-specified ef value as input. Instead, it dynamically determines an appropriate ef during the search process.

Ada-ef extends the standard HNSW search described in Section 3, with several key distinctions elaborated below. First, the input to Ada-ef does not require a predefined ef value; instead, ef is initially set to $\infty$, allowing the best-first graph traversal to visit any nodes around the entry point $ep$. Additionally, the algorithm initializes a sampled distance list $D$ to record the distances between visited nodes and the query vector $\mathbf{q}$. The size of $D$ is bounded by $l$ that is the number of nodes reachable within 2-hops from the entry point $ep$ at the base layer. Once $l$ distances have been collected, the algorithm invokes the ESTIMATE-EF function (Algorithm 1 introduced later), which takes the query vector $\mathbf{q}$, the list $D$, and the user-specified recall $r$ as input, and returns a dynamically estimated ef value. Due to space constraint, we provide the pseudocode of the Ada-ef search algorithm in the appendix.

The ef-estimator, referred to as ESTIMATE-EF, serves as the core component of the Ada-ef search algorithm. We provide a high-level overview of its design, with a detailed presentation in the subsequent sections. The estimator is grounded in a theoretical analysis of the similarity distance distribution. Specifically, we demonstrate that given a query vector $\mathbf{q}$ and a dataset $\mathbf{V}$ it is feasible to approximate the distribution of the full distance list $FDL(\mathbf{q}, \mathbf{V})$ without computing each individual distance. The theoretical basis and methodology for estimating this distribution are detailed in Section 5. Given the estimated $FDL(\mathbf{q}, \mathbf{V})$ distribution, we can infer the $p$-th percentile distance among all the distances. This inference is lightweight and does not require computing any actual distances. Then, we compute the number of candidate distances in the list $D$ that fall below this estimated $p$-th quantile threshold. For instance, if $p = 0.001$, we count how many of the candidate distances lie within the smallest 0.1% of all distances in $FDL(\mathbf{q}, \mathbf{V})$. This count is used to characterize the query difficulty. Specifically, for a given query vector $\mathbf{q}$, a high count, particularly one approaching $|D|$, indicates that many candidate neighbors are highly similar to $\mathbf{q}$, and thus a small ef value is likely sufficient to achieve high recall. Conversely, a low count might suggest that a larger ef is required to reach the same recall level. We leverage this information to estimate an appropriate ef value for each incoming query. The technical details of this process are presented in Section 6.

## 5 FDL Distribution

We analyze the FDL distribution under commonly used similarity metrics for high-dimensional learned embedding vectors, in particular those derived from Transformer-based models, including inner product, cosine similarity, and cosine distance. Our analysis begins with the inner product and is then systematically extended to the others. Our key finding is that, across all considered metrics, $FDL(\mathbf{q}, \mathbf{V})$ approximately follows a Gaussian (normal) distribution when dimensionality becomes large. We derive closed-form expressions for the mean and variance of this distribution. We then discuss how to estimate the distribution efficiently at runtime.

### 5.1 Inner Product

The inner product (IP) between a query vector $\mathbf{q} = (q_1, \ldots, q_d)$ and a data vector $\mathbf{v} = (v_1, \ldots, v_d)$ is defined as: $\mathbf{q} \cdot \mathbf{v} = \sum_{i=1}^{d} q_i v_i$. Given a

dataset $\mathbf{V} = (\mathbf{v}_1, \ldots, \mathbf{v}_n)$, we define the FDL based on inner product as: $FDL_{IP}(\mathbf{q}, \mathbf{V}) = (\mathbf{q} \cdot \mathbf{v}_1, \ldots, \mathbf{q} \cdot \mathbf{v}_n)$.

We begin by introducing a property of $\mathbf{V}$. Under the assumption that this property holds, we then demonstrate that $FDL_{IP}(\mathbf{q}, \mathbf{V})$ follows a normal distribution.

*Definition 5.1.* Let $\mathbf{V}$ be a vector dataset of $n$ data vectors, each of dimensionality $d$. Let $\mathbf{v} = (v_1, \ldots, v_d)$ be a random data vector of $d$ random variables, where each $v_i$ represents the distribution of the values within the $i$-th column of $\mathbf{V}$. If the random variables $(v_1, \ldots, v_d)$ are independent and identically distributed (i.i.d.), then we say the dataset $\mathbf{V}$ is i.i.d across its dimensions.

THEOREM 5.2. *Let $\mathbf{q} = (q_1, \ldots, q_d)$ be a query vector, and let $\mathbf{V}$ be a dataset of data vectors, each of dimensionality $d$, and $\mathbf{V}$ is i.i.d. across its dimensions. Then, the full distance list $FDL_{IP}(\mathbf{q}, \mathbf{V})$ converges in distribution to a normal distribution as $d \to \infty$: $\mathcal{N}\left(\mu_{IP}, \sigma_{IP}^2\right)$ with the mean $\mu_{IP}$ and the variance $\sigma_{IP}^2$: $\mu_{IP} = \sum_{i=1}^{d} q_i E[v_i]$; $\sigma_{IP}^2 = \sum_{i=1}^{d} q_i^2 Var(v_i)$.*

PROOF. Consider the inner product $\mathbf{q} \cdot \mathbf{v} = \sum_{i=1}^{d} q_i v_i$. The query vector $\mathbf{q}$ is given during similarity search online, such that each $q_i$ is a constant at query time. Then, the inner product is a sum of i.i.d random variables, where each is given by $q_i v_i$. According to the Central Limit Theorem, the sum follows a normal distribution as $d$ becomes large. For the mean and variance, we have: $E[\sum_{i=1}^{d} q_i v_i] = \sum_{i=1}^{d} q_i E[v_i]$, and $Var(\sum_{i=1}^{d} q_i v_i) = \sum_{i=1}^{d} q_i^2 Var(v_i)$. □

Theorem 5.2 shows that $FDL_{IP}(\mathbf{q}, \mathbf{V})$ follows a normal distribution under the assumption that the components of $\mathbf{V}$ are i.i.d. across dimensions. This assumption facilitates the application of the classical Central Limit Theorem in the proof of Theorem 5.2. We provide empirical evidence that the identical distribution assumption is approximately valid in real-world datasets. As shown in Figure 3, the dimensions of GloVe and MS MARCO embeddings exhibit distributions that are approximately normal, with minor variations in mean and variance. Similar distributional patterns have been empirically observed in real-world datasets [19, 37, 68], and have also been explicitly encouraged through model design in recent work [58]. Therefore, we adopt the identical distribution assumption in this work. For scenarios where the distributional assumption may not hold, we provide a detailed discussion in the appendix, leveraging Lindeberg's condition [57].

*Relaxing independent condition.* Both the classical CLT and its relaxed form under Lindeberg's condition require that the random variables involved are mutually independent, *i.e.*, each pair $q_i v_i$ and $q_j v_j$ must be independent for $i \neq j$. In practice, however, data vectors are typically learned embeddings, where each dimension represents a latent feature. While these features are ideally independent, in reality, dimensions may exhibit correlations. Such correlations do not affect $\mu_{IP}$, but they do influence $\sigma_{IP}^2$. To account for these correlations, we adjust the variance as follows:

$$\mathrm{Var}\left(\sum_{i=1}^{d} q_i v_i\right) = \sum_{i=1}^{d} q_i^2 \mathrm{Var}(v_i) + 2 \sum_{1 \leq i < j \leq d} q_i q_j \mathrm{Cov}(v_i, v_j).$$
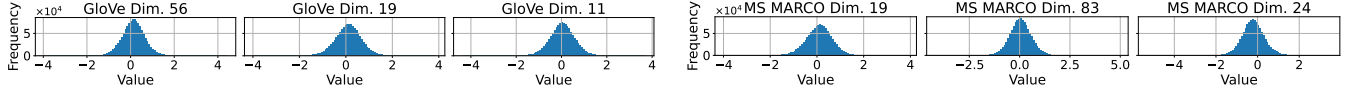
**Figure 3: Distribution of individual embedding dimensions that are randomly sampled in GloVe and MS MARCO.**

Therefore, we model $FDL_{IP}(\mathbf{q}, \mathbf{V})$ as a normal distribution with the adjusted variance:

$$\mathcal{N}\left(\mu_{IP}, \sigma_{IP}^2 + \Delta_{IP}\right), \Delta_{IP} = 2 \sum_{1 \le i < j \le d} q_i q_j \text{Cov}(X_i, X_j). \quad (1)$$

## 5.2 Cosine Similarity

The cosine similarity (CS) between $\mathbf{q}$ and $\mathbf{v}$ is $CS(\mathbf{q}, \mathbf{v}) = \frac{\mathbf{q} \cdot \mathbf{v}}{\|\mathbf{q}\| \|\mathbf{v}\|}$. Given a dataset $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$, the full distance list based on CS is $FDL_{CS}(\mathbf{q}, \mathbf{V}) = (CS(\mathbf{q}, \mathbf{v}_1), \dots, CS(\mathbf{q}, \mathbf{v_n}))$.

We analyze the distribution of $FDL_{CS}(\mathbf{q}, \mathbf{V})$ based on Theorem 5.2. Cosine similarity can be interpreted as the normalized inner product, *i.e.*, $CS(\mathbf{q}, \mathbf{v}) = \hat{\mathbf{q}} \cdot \hat{\mathbf{v}}$, where $\hat{\mathbf{q}} = \frac{\mathbf{q}}{\|\mathbf{q}\|}$ and $\hat{\mathbf{v}} = \frac{\mathbf{v}}{\|\mathbf{v}\|}$. Let $\hat{\mathbf{q}} = (\hat{q}_1, ..., \hat{q}_d)$ and $\hat{\mathbf{v}} = (\hat{v}_1, ..., \hat{v}_d)$, where $\hat{q}_i = \frac{q_i}{\|\mathbf{q}\|}$ and $\hat{v}_i = \frac{v_i}{\|\mathbf{v}\|}$.

According to Theorem 5.2, if the random variables $(\hat{v}_1, \dots, \hat{v}_d)$ representing the normalized column-wise distribution of $\mathbf{V}$ are i.i.d. with finite means and variances, then as dimension $d$ becomes large, the distribution of $FDL_{CS}(\mathbf{q}, \mathbf{V})$ converges to a normal distribution as: $\mathcal{N}\left(\mu_{CS}, \sigma_{CS}^2\right)$, with the following mean and variance:

$$\mu_{CS} = \sum_{i=1}^{d} \hat{q}_i \text{E}[\hat{v}_i]; \sigma_{CS}^2 = \sum_{i=1}^{d} (\hat{q}_i)^2 \text{Var}(\hat{v}_i).$$

Following the same analysis as in the case of IP, *i.e.*, by relaxing the i.i.d. conditions, $FDL_{CS}(\mathbf{q}, \mathbf{V})$ approximately follows:

$$\mathcal{N}\left(\mu_{CS}, \sigma_{CS}^2 + \Delta_{CS}\right), \Delta_{CS} = 2 \sum_{1 \le i < j \le d} \hat{q}_i \hat{q}_j \text{Cov}(\hat{v}_i, \hat{v}_j). \quad (2)$$

## 5.3 Cosine Distance

The cosine distance (CD) between $\mathbf{q}$ and $\mathbf{v}$ is defined as $CD(\mathbf{q}, \mathbf{v}) = 1 - CS(\mathbf{q}, \mathbf{v})$. Accordingly, the full distance list based on CD is defined as: $FDL_{CD}(\mathbf{q}, \mathbf{V}) = (1 - CS(\mathbf{q}, \mathbf{v}_1), \dots, 1 - CS(\mathbf{q}, \mathbf{v_n}))$. Since $FDL_{CS}(\mathbf{q}, \mathbf{V})$ is approximately normally distributed, and since an affine transformation of a normally distributed variable results in another normal distribution (with mean shifted and variance preserved), then $FDL_{CD}(\mathbf{q}, \mathbf{V})$ approximately follows:

$$\mathcal{N}(1 - \mu_{CS}, \sigma_{CS}^2 + \Delta_{CS}). \quad (3)$$

## 5.4 Computing FDL Distribution

Given a query vector $\mathbf{q}$, we aim to efficiently estimate the FDL distribution online by precomputing certain statistics of the dataset $\mathbf{V}$ offline. The core idea is that for the mean and variance of the normal distribution that are formally expressed in Equations (1), (2), and (3), the components depend solely on $\mathbf{V}$ and hence can be computed in advance. At query time, these precomputed statistics are then combined with the query vector $\mathbf{q}$ to enable on-the-fly estimation of the FDL distribution. To illustrate the offline computation of $\mathbf{V}$'s statistics, we focus on the case of the inner product.

*Offline computation.* According to Equation (1), the estimation of $FDL_{IP}$ requires the mean $\mu_{IP}$ and the combined variance term $\sigma_{IP}^2 + \Delta_{IP}$. To facilitate efficient computation during query time, we precompute the following statistics of the dataset $\mathbf{V}$ in the offline phase. Let $v_i$ denote the $i$-th column of $\mathbf{V}$.

- Mean Vector of $\mathbf{V}$: A vector of dimension $1 \times d$, where the $i$-th entry corresponds to $\text{E}[v_i]$, *i.e.*, the mean of the $i$-th column of $\mathbf{V}$.
- Covariance Matrix of $\mathbf{V}$: A matrix of dimension $d \times d$, where the $(i, j)$-th entry represents $\text{Cov}(v_i, v_j)$, *i.e.*, the covariance between the $i$-th and $j$-th columns of $\mathbf{V}$. The diagonal entries of this matrix correspond to $\text{Var}(v_i)$, *i.e.*, the variance of the $i$-th column.

Figure 2 illustrates the mean vector and the covariance matrix.

The mean vector is computed by scanning each column of $\mathbf{V}$ and calculating its mean. To compute the covariance matrix, we proceed as follows. Let $\mathbf{M}$ be a matrix where each row is a copy of the mean vector, *i.e.*, $\mathbf{M}$ replicates the mean vector across all $n$ rows. The covariance matrix $\Sigma$ of $\mathbf{V}$ is then given by: $\Sigma = \frac{1}{n-1} (\mathbf{V} - \mathbf{M})^{\top} (\mathbf{V} - \mathbf{M})$.

*Online computation.* Given a query vector $\mathbf{q}$, in order to estimate $FDL_{IP}(\mathbf{q}, \mathbf{V})$, we use the mean vector and covariance matrix of the dataset $\mathbf{V}$ to compute $\mu_{IP}$ and $\sigma_{IP}^2 + \Delta_{IP}$, as defined in Theorem 5.2 and Equation 1. We compute $\mu_{IP}$ as the dot product between the query $\mathbf{q}$ and the mean of $\mathbf{V}$, and compute $\sigma_{IP}^2 + \Delta_{IP}$ as $\mathbf{q}\Sigma\mathbf{q}^{\top}$.

Although focused on the inner product, the approach extends to cosine similarity and distance. Both offline and online computations can be formulated as matrix operations, enabling efficient SIMD and parallel acceleration.

## 6 Query-Aware EF Estimation

We present the design of the ef-estimator, *i.e.*, the ESTIMATE-EF function. The primary objective of the ef-estimator is to *differentiate between query vectors that can achieve the declarative target recall with a small ef value and those that require a large ef value to obtain the recall*. To accomplish this, we utilize the estimated distribution of the FDL to design a query scoring model, *i.e.*, a mechanism that assigns a score to each input query vector. We estimate ef for each input query based on its query score. To achieve this, we perform offline sampling of data vectors from the dataset $\mathbf{V}$ to build an ef-estimation table. This table stores pairs of ef values and their corresponding estimated recall values, indexed by query scores. During the online search phase, the query scoring mechanism is applied to compute the score for the incoming query, and then the appropriate ef value is retrieved from the ef-estimation table based on this score. In essence, the query scoring mechanism distinguishes queries requiring different ef values to meet the target recall, while the ef-estimation table provides the actual ef values needed to achieve it. We present the query scoring mechanism in

§ 6.1, the ef-estimation table in § 6.2, and discuss how to perform incremental insertion and deletion in § 6.3.

## 6.1 Query Scoring Model

The adaptive-ef search visits a fixed number of nodes at the base layer of the HNSW index and collects a distance list $D$, which contains the distances between the query vector $\mathbf{q}$ and the visited nodes, as described in § 4. The query scoring model analyzes the list $D$ to assess the difficulty of computing the approximate nearest neighbors for $\mathbf{q}$. The design of the scoring model is guided by the following principle: *if the majority of distances in $D$ fall among the smallest distances according to the distribution of $FDL(\mathbf{q}, \mathbf{V})$, then $\mathbf{q}$ likely requires only a small ef value to achieve high recall. Otherwise, $\mathbf{q}$ may require a larger ef value for specified recall.*

The query scoring model discretizes the normal distribution (parameterized by the estimated mean and variance of the FDL) into a fixed number of consecutive quantile-based bins, each covering an equal portion of the FDL referred to as $\delta$. These bins are denoted as $(b_1, \ldots, b_m)$, where $m$ is the total number of bins. Each bin is used to count the number of distances that fall within the quantile range $\delta$ of the distribution. For instance, if $\delta = 0.001$, $b_1$ corresponds to distances smaller than the 0.001 quantile of the FDL distribution, while $b_2$ captures distances between the 0.001 and 0.002 quantiles, and so on for the remaining bins.

The bins $(b_1, \ldots, b_m)$ are formally defined by threshold values derived from the quantiles of the FDL distribution. The threshold $\theta_i$ for separating $b_i$ and $b_{i+1}$ is computed as:

$$\theta_i = \mu + \sigma\Phi^{-1}(\delta i), \text{ for } i = 1, \ldots, m, \tag{4}$$

where $\Phi^{-1}(\cdot)$ denotes the inverse cumulative distribution function, $\delta$ multiplied by $i$ determines the quantile, and $\mu$ and $\sigma$ are the mean and the standard deviation of the normal distribution, respectively, *e.g.*, for cosine distance, $\mu = 1 - \mu_{CS}$ and $\sigma = \sqrt{\sigma_{CS}^2 + \Delta_{CS}}$. Each bin $b_i$ is then associated with a weight $w_i$ to indicate the importance of the bin, which is defined via an exponential decay function: $w_i = 100e^{-i+1}$, for $i = 1, \ldots, m$.

Given a list $D$ containing the distances between the query vector and the visited nodes, we count how many falling into each bin:

$$c_i = \left|\left\{\theta_{i-1} < d_j \leq \theta_i \mid d_j \in D\right\}\right|. \tag{5}$$

The query score $s(\mathbf{q})$ is a weighted, normalized sum of bin counts:

$$s(\mathbf{q}) = \sum_{i=1}^{m} \left(w_i \frac{c_i}{|D|}\right). \tag{6}$$

This query score serves as a proxy for query difficulty: a higher score indicates that more distances in $D$ fall within the favorable (low) quantile regions of the FDL distribution. Thus, a small ef would be sufficient for high-score queries to achieve high recall. We provide an example to explain the query scoring model in detail in the appendix.

We note that several query-difficulty metrics have been proposed, including LID [16], RC [48], Expansion [15], and Steiner-Hardness [93], which help characterize query hardness, but are impractical for online search. They require exact neighbor information or fine-grained distance statistics obtained through additional

nearest-neighbor or multi-radius searches. In contrast, our FDL-based approach is lightweight and can estimate query difficulty directly during the online search process.

## 6.2 Estimate EF

The query scoring model serves to quantitatively differentiate between queries that require a small ef value and those that require a larger one. However, the model alone does not directly provide a concrete ef value for a given query at runtime. To bridge this gap, we construct an ef-estimation table offline, which records the relationship between query scores, ef values, and the corresponding recall levels. At query time, we compute the score for each incoming query vector and use it to retrieve from the ef-estimation table an ef value that is expected to achieve the target recall.

*Offline ef-estimation table construction.* A key challenge in constructing the ef-estimation table is the absence of query vectors during the offline phase. To address this, we uniformly sample a subset of data vectors (*e.g.*, 200 vectors) from the dataset $\mathbf{V}$ and use them as proxy query vectors. For each sampled vector, we compute its query score and group the vectors according to their scores. To reduce the granularity and number of score groups, the floating-point query scores are casted into integer ones. For each score group, we evaluate recall using the original HNSW search under progressively increasing ef values. If the observed average recall for a group falls below the user-defined target, a larger ef value is selected and tested. This adaptive probing continues until the target recall is achieved or a predefined upper bound on ef is reached. The resulting mappings between query score, ef, and recall are stored in the ef-estimation table.

Figure 2 shows an example of ef-estimation table. Specifically, the table consists of two columns: Score and EF_Recall. The Score column represents the query score group, while the EF_Recall column contains a list of (EF, Recall) pairs, sorted in ascending order of ef. Each entry in this list captures the recall achieved for a specific ef value within the corresponding score group.

*Weighted Average ef.* We compute the weighted average ef (WAE) as a summary of ef-estimation table, weighted by the number of queries in each score group (*i.e.*, an integer query score): WAE $= \frac{1}{G}\sum_i g_i \cdot \text{ef}_i$, where $g_i$ is the number of queries in score group $i$, $\text{ef}_i$ is the smallest *ef* that achieves the target recall for that group, and $G$ is the total number of sampled data vectors. This metric captures the average search effort required across queries of varying difficulty to reach the target recall and serves as an initialization point for ef estimation.

*Online ef estimation.* Algorithm 1 presents the end-to-end procedure for estimating ef for a given query vector $\mathbf{q}$ to achieve a target recall $r$. The process begins by estimating the parameters of the distribution $FDL(\mathbf{q}, \mathbf{V})$, specifically the mean $\mu$ and standard deviation $\sigma$ (lines 1–2). These computations are grounded in the theoretical foundations described in § 5.4, including Equations (1), (2), and (3). For efficient on-the-fly computation, the algorithm utilizes the mean vector and covariance matrix of the dataset $\mathbf{V}$, both of which are precomputed during the offline phase. Lines 3–5 focus on computing the query score based on the input distance list $D$ and the estimated parameters $\mu$ and $\sigma$. The algorithm first constructs a set of quantile-based bins $(b_1, \ldots, b_m)$ (line 3) and then determines

**Algorithm 1** ESTIMATE-EF($\mathbf{q}, D, r$)

---

**Input:** query vector $\mathbf{q}$, distance list $D$, expected recall $r$
**Output:** estimated exploration factor ef

1: $\mu \leftarrow$ compute the estimated mean of $FDL(\mathbf{q}, \mathbf{V})$      ▷ § 5.4
2: $\sigma \leftarrow$ compute the estimated std of $FDL(\mathbf{q}, \mathbf{V})$      ▷ § 5.4
3: $(b_1, \ldots, b_m) \leftarrow$ compute the bins      ▷ Eq. (4)
4: $(c_1, \ldots, c_m) \leftarrow$ compute the bin counts      ▷ Eq. (5)
5: $s \leftarrow$ compute the query score      ▷ Eq. (6)
6: $row \leftarrow$ get the row in ef-estimation table with the score
7: ef $\leftarrow$ get the largest EF in the $row$
8: **for** each (EF, Recall) in the $row$ **do**
9:      **if** Recall $\geq r$ **then**
10:          ef $\leftarrow$ max(EF, WAE)
11:          **break**
12: **return** ef

---

the bin counts $(c_1, \ldots, c_m)$ by analyzing how the distances in $D$ fall into these bins (line 4). These counts are subsequently aggregated to compute a scalar query score $s$ (line 5), which serves as a proxy for the difficulty of the query. After computing the score, the algorithm retrieves the corresponding row from the precomputed ef-estimation table, which associates each score group with a list of ($EF, Recall$) pairs sorted in ascending order of ef (lines 6–7). It then searches this list to find the smallest ef value that satisfies the target recall requirement (lines 8–11). If the selected value is smaller than WAE, it is increased to WAE and returned as the estimated ef; otherwise, if no such value exists, the algorithm returns the largest ef available for the specific score.

We note the query-aware ef estimation framework is built on the FDL distribution. We have derived FDL distributions for Inner Product, Cosine Similarity, and Cosine Distance, while the Euclidean (L2) case remains open due to the squared terms in its formulation. Once established, the framework could be applied to L2 as well.

## 6.3 Incremental Updates

In practice, vector indexes such as HNSW are expected to support dynamic vector insertion and deletion, $i.e.$, the ability to insert into and delete from the index. We outline how our approach accommodates such updates. Offline computation in our method consists of three main components: (i) dataset-level statistics, including the mean vector and covariance matrix (see Section 5.4); (ii) a sampled subset of data vectors and computing their ground truth; (iii) the ef-estimation table. To support vector updates, we observe that components (i) and (ii) can be incrementally updated. Once these are updated, component (iii), $i.e.$, the ef-estimation table can be recomputed with respect to the updated HNSW index. Updating the ground truth of sampled vectors is straightforward. For insertions, we compute distances to newly added vectors and update the Top-$k$ neighbors. For deletions, we incrementally remove deleted vectors from the ground truth and refresh the Top-$k$ neighbors accordingly. We now briefly describe how to incrementally update the dataset statistics, beginning with the case of insertions. Let $\mathbf{V}$ denote the original dataset of $n$ vectors, and let $\mathbf{V}'$ represent the set of $n'$ newly inserted vectors. The updated dataset is then given by $\mathbf{V}'' = \mathbf{V} \cup \mathbf{V}'$,

**Table 1: Overview of real-world datasets.**

| Dataset | Data Size | Dimension | Query Size |
|---|---|---|---|
| GloVe-100 [75] | 1,183,514 | 100 | 10,000 |
| DeepImage [94] | 9,990,000 | 96 | 10,000 |
| MS MARCO V1 [11] | 8,841,823 | 1536 | 6980 |
| MS MARCO V2.1 [14] | 18,380,565 | 1024 | 1677 |
| Laion-I2I [80] | 30,646,115 | 512 | 10,000 |
| Laion-T2I [80] | 30,646,115 | 512 | 10,000 |

containing $n'' = n + n'$ vectors. We first compute the mean vector $\mathbf{M}'$ and covariance matrix $\Sigma'$ of $\mathbf{V}'$, and then merge them with the existing $\mathbf{M}$ and $\Sigma$ of $\mathbf{V}$ as follows: $\mathbf{M}'' = \frac{n\mathbf{M} + n'\mathbf{M}'}{n''}$, and $\Sigma'' = \frac{1}{n''-1} \left[ (n-1)\Sigma + (n'-1)\Sigma' + \frac{nn'}{n''}(\mathbf{M} - \mathbf{M}')^\top (\mathbf{M} - \mathbf{M}') \right]$. Consider deleting $\mathbf{V}'$ from $\mathbf{V}''$ to obtain $\mathbf{M}$ and $\Sigma$ of $\mathbf{V}$: $\mathbf{M} = \frac{n''\mathbf{M}'' - n'\mathbf{M}'}{n}$, and $\Sigma = \frac{1}{n-1} \left[ (n''-1)\Sigma'' - (n'-1)\Sigma' - \frac{n'n''}{n}(\mathbf{M}'' - \mathbf{M}')^\top (\mathbf{M}'' - \mathbf{M}') \right]$. We note that such online updates to statistical summaries are widely adopted in practice, $e.g.$, in Spark MLlib [66].

## 7 Experimental Evaluation

We evaluate online effectiveness (§ 7.2) and offline computation (§ 7.3) using six real-world and two synthetic datasets against five baselines, including state-of-the-art adaptive methods. We report sensitivity analysis (§ 7.4), update performance (§ 7.5), and ablation study (§ 7.6).

## 7.1 Experiment Setup

*Real-world datasets.* We use six real-world datasets, summarized in Table 1. From the early ANNS benchmark [20], we include GloVe-100 (word) [75] and DeepImage (image) [94]. To reflect more recent trends, we incorporate several modern datasets that use advanced embedding techniques. Specifically, we include the MS MARCO V1 Passage Ranking dataset [24], a widely used benchmark for passage reranking, which is embedded using OpenAI's ada2 [46]. This embedding captures semantic relationships between passages and queries and is particularly suited for dense retrieval tasks. This version is sourced from the Anserini framework [11, 96]. We also use MS MARCO V2.1, drawn from the TREC-RAG 2024 corpus [13], embedded using Cohere's Embed English V3. This model produces dense semantic embeddings of passages, optimized for retrieval-augmented generation (RAG) settings and downstream QA applications. Finally, we evaluate our approach in a multi-modal retrieval setting using the LAION dataset [80], which consists of text-image pairs embedded with OpenAI's CLIP [78]. It aligns visual and textual modalities in a shared embedding space, enabling semantic similarity search across modalities. We consider two retrieval tasks: image-to-image retrieval and text-to-image retrieval, referred to as Laion-I2I and Laion-T2I, respectively. Both tasks utilize the same set of image embeddings as the database, but differ in the query modality: Laion-I2I uses image queries, while Laion-T2I uses text queries. Each dataset comes with a set of query vectors, except for LAION for which we randomly sample 10000 image embedding vectors and 10000 text embedding vectors as query vectors, respectively.

*Synthetic datasets.* To evaluate the effectiveness of our adaptive approach, we include two synthetic datasets, Uniform Cluster and Zipfian Cluster. Both datasets contain 10 million vectors in 100

dimensions, organized into 5000 Gaussian distributed clusters. The primary difference between them lies in the distribution of cluster sizes. In Uniform Cluster, all clusters have the same size, with each containing exactly 2000 vectors. In contrast, Zipfian Cluster follows a Zipfian distribution with an exponent of 1, resulting in highly skewed cluster sizes. The largest cluster contains more than one million vectors, and the sizes decrease exponentially. For each dataset, we sample 10000 vectors to serve as query vectors for evaluating retrieval performance.

*Implementation.* We integrate our approach into `HNSWlib` [2], the official implementation of HNSW developed by its original authors. We use Eigen [3] for efficient matrix operations involved in dataset statistics computation. Our method, referred to as **Adaptive-ef** (**Ada-ef**), extends the HNSW index without modification. Our source code[1] is publicly available online.

*Baselines.* We include the following approaches as baselines: Learned Adaptive Early Termination (**LAET**) [59], **DARTH** [27], Patience in Proximity (**PiP**) [82], **HNSW (HNSWlib)** [2], and **HNSW (FAISS)** [1]. Both LAET and DARTH are learning-based adaptive vector search methods built on Gradient Boosting Decision Trees (GBDTs) [39]. They share similar feature designs and both require offline preparation of training data and model training. The key difference is that DARTH periodically predicts the current recall during search to determine whether the declarative target recall is met, whereas LAET performs a single prediction to estimate the required number of distance computations. We use the open-source implementation of DARTH, which also provides a declarative recall variant of LAET. PiP is an early-termination method that requires no offline computation, which terminates the search if the current results have not been updated for a fixed amount of consecutive steps. We implement it on top of `HNSWlib`. Finally, we include two widely used HNSW implementations, HNSW (HNSWlib) and HNSW (FAISS), to evaluate the benefits of adaptive approaches. We note that DARTH represents the state of the art among these baselines [27].

*Parameters.* For the construction of HNSW indexes on each dataset, we set $M = 16$ (maximum number of outgoing degree) and $efConstruction = 500$ (search parameter during indexing). For Top-$k$ ANNS, we follow standard settings in existing benchmarks. Specifically, we set $k = 1000$ for both the MS MARCO and LAION datasets, following recommendations from the MS MARCO benchmark [7] and the SISAP Indexing Challenge [12]. This setting mirrors typical retrieval pipelines in modern applications, such as RAG systems [44], where a recall stage first retrieves a broad set of candidate items via Top-$k$ ANNS, followed by a reranking stage that refines these candidates for improved relevance. For the remaining datasets from the ANNS benchmark suite [4], we adopt the commonly used setting of $k = 100$. The influence of varying $k$ is analyzed in Section 7.4. We use cosine distance in all experiments, as it is a widely adopted choice for ANNS over high-dimensional Transformer-based learned embeddings. For DARTH and LAET, we train the models using 10000 learn vectors, following the recommended configuration [27]. The prediction intervals in DARTH and the number of distance computations at which LAET collects features for training or predicting follow the same setup as in DARTH, *i.e.*, extracting the statistics of the training data. For PiP, we set

---

[1] **GitHub Repository:** https://github.com/chaozhang-cs/hnsw-ada-ef

the saturation threshold to 0.95, as recommended in the original paper, and the patience parameter to 30. Since the configuration of patience was not explicitly specified in PiP [82], we adopt a conservative value to ensure stable convergence. We note that Ada-ef requires no parameter setting, except in the ablation studies in § 7.6, which examine its design choices.

*Settings.* All experiments are conducted on a shared server running Ubuntu 22.04, equipped with 80 physical CPU cores (Intel Xeon Platinum 8380 @ 2.30 GHz), 160 hardware threads, and 1 TB of DDR4-3200 MHz RDIMM main memory. All implementations are written in C++, leveraging SIMD intrinsics for vectorized computations. The code is compiled with g++ 12.3.0, using the optimization flag -O3 and enabling AVX-512, via setting `FAISS_OPT_LEVEL` as avx512 for FAISS-based methods—HNSW (FAISS), DARTH, and LAET—or setting -march=native for HNSWlib-based ones—HNSW (HNSWlib), Ada-ef, and PiP. All offline experiments are conducted using 40 threads. For online searching, we follow the common practice of using a single thread [27]. Note that queries can be processed in inter-query parallelism to fully utilize all available CPU cores.

## 7.2 Online Searching Performance

We present the search results in Figure 4, which reports both recall and end-to-end query workload execution time. For DARTH, LAET, and Ada-ef, we set the target recall to 0.95 across all datasets. For PiP, we set ef to 300 for Top-100 ANNS or 2000 for Top-1000 ANNS. For HNSW (HNSWlib) and HNSW (FAISS), we execute the query workload multiple times, gradually increasing ef starting from ef = $k$ for Top-$k$ ANNS, until either the average recall reaches 0.99 or ef reaches a upper bound of 5000, and we report the execution time and corresponding recall. We note that this experiment focuses on simulating a real-time scenario where queries are served on the fly, and HNSW (HNSWlib) and HNSW (FAISS) results under varying ef values are used as references because tuning ef values requires ground-truth for incoming queries. In Figure 4, we also highlight key reference points to aid interpretation. Specifically, we mark the cases of ef = $k$ and ef = $2k$ using two vertical dotted gray lines; the former is one of the commonly used settings in practice. The target recall of 0.95 is indicated with a horizontal dashed line. To better assess search quality, we report three recall statistics: average recall, 1st percentile, and 5th percentile. While average recall reflects the overall workload performance, the 1st and 5th percentiles capture performance on the hardest queries. These metrics help diagnose under-searching and over-searching behaviors. Figure 5 shows the distribution of ef values dynamically configured by Ada-ef.

For HNSW (HNSWlib) and HNSW (FAISS) results, different datasets require different ef values to achieve high recall, and a fixed setting of ef = $k$ is insufficient. Moreover, increasing ef beyond a certain threshold results in diminishing returns, *i.e.*, recall improves only marginally while query time increases substantially, highlighting the over-searching issue.

For PiP, the achieved recall is substantially lower than that of the other methods, highlighting the limitation of its fixed-threshold early termination strategy in maintaining practical effectiveness.

For the learning-based adaptive approaches, DARTH and LAET, both methods reach the declarative target recall of 0.95 on five out of the eight evaluated datasets. However, on these datasets, they
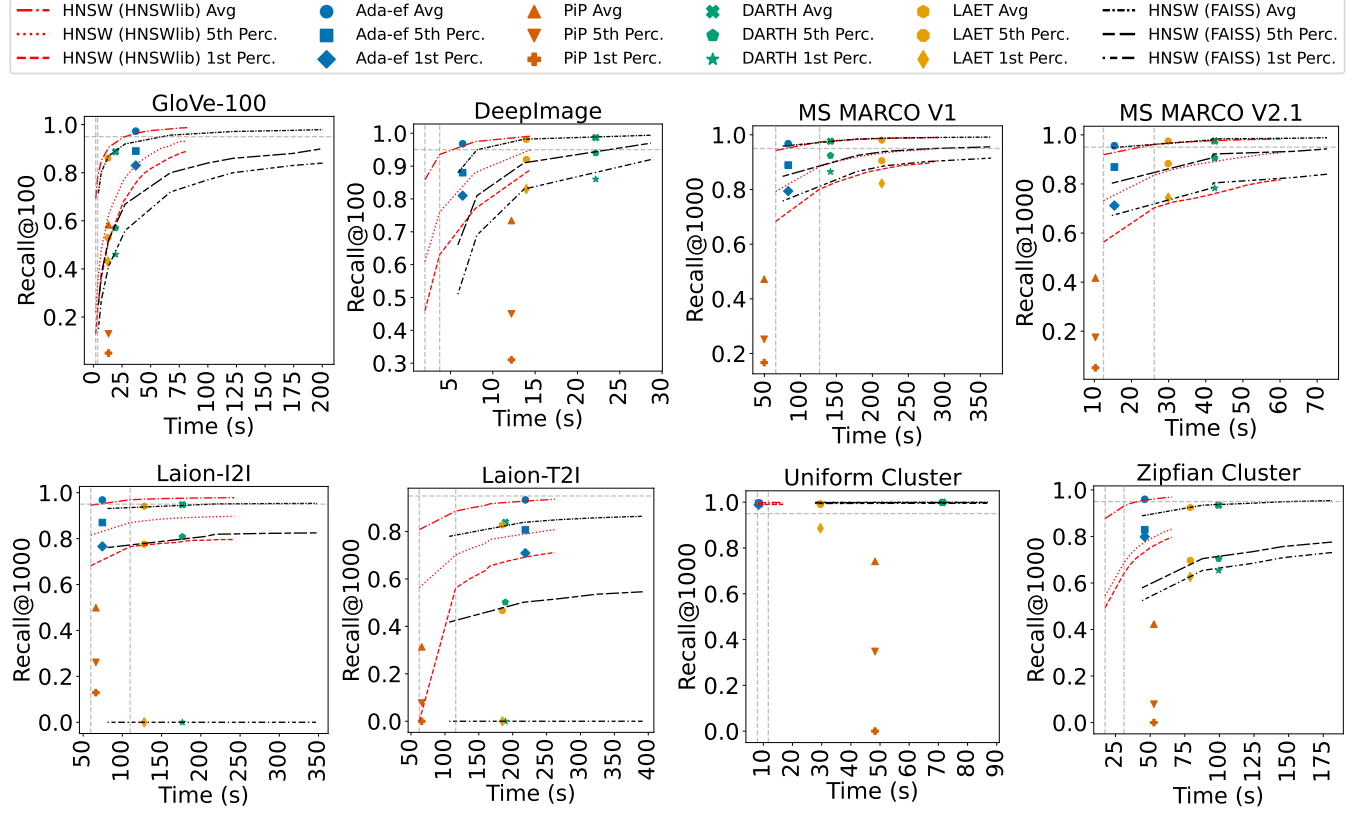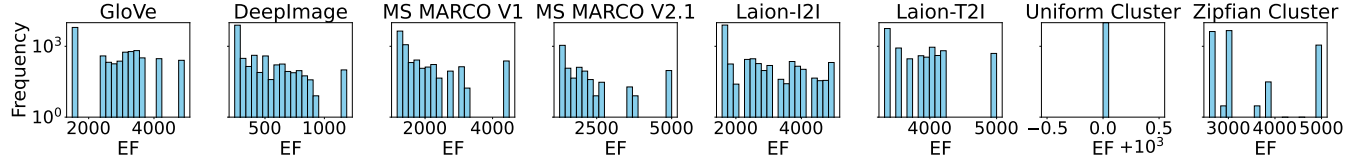
Figure 4: Search performance on real and synthetic datasets.



Figure 5: Distribution of `ef` values dynamically assigned by `Ada-ef` across queries in each dataset (log scale).

tend to over-search, entering a plateau region where additional computation produced negligible recall gains, *e.g.*, MS MARCO V1. For the remaining datasets, their recall fails to meet the target and is, in some cases, considerably lower, *e.g.*, GloVe-100. Overall, DARTH generally achieves higher recall than LAET, but at the expense of substantially higher query latency.

`Ada-ef` consistently achieves the best performance across all datasets. It dynamically configures the `ef` value on a per-query basis to approximately meet the declarative target recall, effectively avoiding over-searching by skipping unnecessarily large `ef` values that offer only marginal recall improvements at the cost of significantly increased latency. At the same time, it mitigates under-searching by improving recall for harder queries, as reflected in the enhanced 1st and 5th percentile recall metrics, which are substantially higher than those of HNSW (HNSWlib). Compared to the state-of-the-art approach DARTH, `Ada-ef` achieves up to a

4× reduction in workload latency for the same level of average recall. In addition, the distributions of dynamically configured `ef` values for each query, shown in Figure 5, exhibit a long-tail pattern, *i.e.*, smaller `ef` values are assigned to the majority of queries, while a minority receive larger values. This fine-grained dynamic configuration substantially improves the worst-case recall.

In the multimodal setting, *i.e.*, Laion-I2I and Laion-T2I, although the dataset vectors (image embeddings) are identical, the differing modalities of the query vectors lead to substantial variations in search performance. The text-to-image scenario is considerably more challenging than the image-to-image case: setting `ef` to $k$ yields a recall close to 0.95 for Laion-I2I, but only around 0.8 for Laion-T2I. This discrepancy arises primarily from the latent distribution gap between visual and textual modalities, which makes cross-modal retrieval inherently more difficult. Existing approaches achieve recall levels well below the declarative target
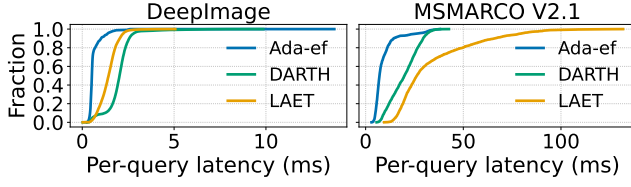
**Figure 6: CDFs of per-query latencies.**

of 0.95, whereas `Ada-ef` adaptively adjusts `ef` based on the joint distribution of data and query vectors, attaining recall close to the target without incurring excessive computation.

The two synthetic datasets, Uniform Cluster and Zipfian Cluster, share identical characteristics except for their cluster size distributions: Uniform Cluster consists of equally sized clusters, whereas Zipfian Cluster exhibits highly skewed cluster sizes following a Zipfian distribution. This difference leads to distinct search behaviors. Achieving high recall is substantially easier in the Uniform case than in the Zipfian case. All existing methods struggle under the Zipfian setting, attaining lower recall while incurring significantly higher search time. In contrast, `Ada-ef` demonstrates robust performance, substantially improving both the 5th and 1st percentile recall and achieving the target recall without unnecessary computation. These results highlight the practical effectiveness of `Ada-ef`, particularly in real-world scenarios where data and query distributions are often skewed rather than uniform.

We use DeepImage and MS MARCO V2.1 as representatives to analyze the per-query latency distribution. The cumulative distribution functions (CDFs) are shown in Figure 6. As illustrated, `Ada-ef` achieves low latency for the majority of queries while spending more time on a small subset of difficult ones. This behavior indicates that `Ada-ef` effectively avoids both over-searching and under-searching, thereby reaching the target recall more efficiently and improving recall performance for the hardest queries.

### 7.3 Offline Computing Performance

Table 2 reports the offline computation time. The `HNSW (HNSWlib)` and `HNSW (FAISS)` columns show the index construction time. For `DARTH` and `LAET`, the offline computation consists of three steps: (i) computing the ground truth for sampled 10000 learn vectors, (ii) generating training data using these vectors, and (iii) training the model with the generated data. These steps are denoted as `LVec GT`, `TData`, and `Train`, respectively. The offline computation of `Ada-ef` includes: (i) computing dataset-level statistics (`Stats`), (ii) sampling 200 data vectors and computing their ground truth (`Samp.`), and (iii) constructing the `ef`-estimation table (`EF-Est.`) with an expected recall of 0.95 and an `ef` upper bound of 5000. Table 3 presents the corresponding memory usage. The `HNSW (HNSWlib)` and `HNSW (FAISS)` columns indicate index size; the `DARTH` and `LAET` columns report the memory footprint of their data and models; and the `Ada-ef` columns show the total memory usage across its components. Note that `PiP` does not require offline computation and is therefore omitted from this experiment.

For the learning-based methods, `DARTH` and `LAET`, the major offline computation overhead is computing the ground-truth for 10000

learning vectors (`LVec GT`), which can exceed the index construction time. In addition, the training data generation stage (`TData`) also incurs a substantial cost. We note that the training time of `LAET` is considerably lower than that of `DARTH`, since `LAET` extracts features once per learn vector, whereas `DARTH` performs this process multiple times for each learn vector.

In contrast, the offline computation time of `Ada-ef` is negligible compared to HNSW index construction, typically less than 5% of the indexing time. Compared with `DARTH` and `LAET`, `Ada-ef` achieves, on average, a 50× reduction in total offline computation time. Its total time is even lower than the least costly offline stage of `DARTH` (model training). These results demonstrate that `Ada-ef` can operate as a lightweight add-on to existing HNSW indexes, providing distribution-aware, adaptive `ef` configuration that improves recall while avoiding both over- and under-searching.

Consider the breakdown of `Ada-ef`'s offline computation cost. We observe that while `Samp.` and `Stats` dominate the offline computation, the cost of constructing the `EF-Est.` table remains negligible. This design highlights the generality and reusability of `Ada-ef`. Specifically, both `Samp.` and `Stats` are agnostic to the value of $k$ and the declarative target recall, meaning that once computed, they can be reused across different configurations, such as in the experiments of Section 7.4, where only the `EF-Est.` table needs to be recomputed, making the additional computation negligible. In contrast, for `DARTH` and `LAET`, both the `TData` generation and model training must be repeated for different $k$ values, resulting in significant offline overhead. Furthermore, both `Stats` and `Samp.` can be efficiently updated incrementally when the dataset is updated (§ 7.5), demonstrating that `Ada-ef` is well-suited for practical, evolving real-world scenarios.

Table 3 reports the offline memory usage with detailed breakdowns for each method. For learning-based approaches, the training data (`TData`) dominates the overall memory footprint. This component primarily consists of large tabular datasets containing the features used to train the GBDT models, which can reach hundreds of millions of rows. In contrast, `Ada-ef` requires minimal offline memory, achieving an average 100× reduction in memory footprint.

### 7.4 Sensitivity Analysis

We analyze the performance of `Ada-ef` under varying configurations, including different values of $k$ for Top-$k$ ANNS and different target recall levels. Specifically, we vary $k$ from 100 to 1000, and for each value, we run `Ada-ef` with three target recall settings: 0.95, 0.97, and 0.99. For comparison, HNSW is executed with increasing `ef` values until either the average recall reaches 0.99 or the maximum `ef` of 5000 is reached. We conduct experiments on two representative datasets: DeepImage and MS MARCO V2.1, and we present the sensitivity analysis results in Figure 7.

Overall, increasing the target recall in `Ada-ef` consistently improves average recall as well as the 5th and 1st percentile recall across different $k$ values and datasets. `Ada-ef` demonstrates robust and reliable behavior, achieving or remaining close to the specified target recall while mitigating both over- and under-searching. Notably, on the DeepImage dataset with $k = 1000$, `Ada-ef` with a target recall of 0.99 incurs higher query latency. However, this additional computation is adaptively allocated to the most difficult

**Table 2: Offline computation time (seconds) with breakdowns for Ada-ef, DARTH, and LAET.**

| Dataset | HNSW (HNSWlib) | Ada-ef | | | | HNSW (FAISS) | DARTH | | | | LAET | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Stats | Samp. | EF-Est. | Total | | LVec GT | TData | Train | Total | LVec GT | TData | Train | Total |
| DeepImage | 439.536 | 0.098 | 6.905 | 0.678 | 7.681 | 602.636 | 179.465 | 37.521 | 157.235 | 374.221 | 179.465 | 37.521 | 0.424 | 217.410 |
| GloVe-100 | 49.238 | 0.014 | 0.794 | 2.125 | 2.933 | 70.937 | 22.820 | 40.618 | 147.395 | 210.833 | 22.820 | 40.618 | 0.439 | 63.877 |
| MS MARCO V1 | 2013.632 | 25.271 | 54.472 | 6.327 | 86.070 | 3148.286 | 2894.502 | 1738.706 | 284.351 | 4917.558 | 2894.502 | 1738.706 | 2.348 | 4635.556 |
| MS MARCO V2.1 | 2263.155 | 13.814 | 80.031 | 8.206 | 102.051 | 4619.897 | 4699.835 | 1819.703 | 381.704 | 6901.242 | 4699.835 | 1819.703 | 1.909 | 6521.447 |
| Laion-I2I | 4110.403 | 4.625 | 71.228 | 3.199 | 79.052 | 5155.223 | 3076.920 | 1793.698 | 226.168 | 5096.786 | 3076.920 | 1793.698 | 0.833 | 4871.451 |
| Laion-T2I | 4110.403 | 4.625 | 69.048 | 6.021 | 79.694 | 5155.223 | 3307.787 | 1688.520 | 274.387 | 5270.693 | 3307.787 | 1688.520 | 0.866 | 4997.173 |
| Uniform Cluster | 198.340 | 0.132 | 7.217 | 0.402 | 7.751 | 241.179 | 193.547 | 141.618 | 32.157 | 367.321 | 193.547 | 141.618 | 0.280 | 335.445 |
| Zipfian Cluster | 436.964 | 0.235 | 7.183 | 2.860 | 10.278 | 522.806 | 210.620 | 1660.215 | 284.941 | 2155.776 | 210.620 | 1660.215 | 0.325 | 1871.160 |

**Table 3: Offline computation memory usage (bytes) with breakdowns for Ada-ef, DARTH, and LAET.**

| Dataset | HNSW (HNSWlib) | Ada-ef | | | | HNSW (FAISS) | DARTH | | | | LAET | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Stats | Samp. | EF-Est. | Total | | LVec GT | TData | Train | Total | LVec GT | TData | Train | Total |
| DeepImage | 5G | 37K | 154K | 2K | 0.2M | 5G | 7.8M | 2G | 342K | 2G | 7.8M | 2G | 294K | 2G |
| GloVe-100 | 620M | 40K | 157K | 3.5K | 0.21M | 0.615G | 7.8M | 2.4G | 348K | 2.4G | 7.8M | 2.4G | 296K | 2.4G |
| MS MARCO V1 | 52G | 9.1M | 2M | 1K | 11.2M | 52G | 98M | 4.7G | 346K | 4.8G | 98M | 4.7G | 397K | 4.8G |
| MS MARCO V2.1 | 73G | 4.1M | 1.6M | 1K | 5.8M | 73G | 79M | 5.2G | 345K | 5.3G | 79M | 5.2G | 360K | 5.3G |
| Laion-I2I | 63G | 1.1M | 1.2M | 2.2K | 3.4M | 63G | 59M | 5G | 342K | 5G | 59M | 5G | 325K | 5G |
| Laion-T2I | 63G | 1.1M | 1.2M | 2.1K | 3.4M | 63G | 59M | 5.6G | 345K | 5.7G | 59M | 5.6G | 321K | 5.7G |
| Uniform Cluster | 5.2G | 40K | 860K | 0.6K | 0.9M | 5.1G | 42.9M | 373M | 323K | 416M | 42.9M | 373M | 280K | 416M |
| Zipfian Cluster | 5.2G | 40K | 860K | 1K | 0.9M | 5.1G | 42.9M | 4.3G | 346K | 4.4G | 42.9M | 4.3G | 295K | 4.4G |



**Figure 7: Sensitivity analysis with different values of $k$ for Top-k ANNS.**

queries, as evidenced by the substantial improvements in the 1st and 5th percentile recall metrics.

## 7.5 Incremental Insertion and Deletion

We use DeepImage and MS MARCO V2.1 as the representative datasets to simulate a practical scenario in which data vectors are inserted into or deleted from the current dataset. Specifically, each original dataset is partitioned into two subsets: the first serves as existing data, and the remaining as updated data, whose size denoted as the batch size (BS) simulates real-world batch updates. We consider two batch sizes: 10% and 50% of the original data.

In the insertion experiments, the update data is inserted into the existing data to reconstruct the original dataset. In this setting, we first build an HNSW index over the existing data and then insert the new data to measure the index update time. For the deletion experiments, we remove the update data from the original dataset to obtain the existing subset. Since HNSW (HNSWlib) does not support in-place deletions (the same holds for HNSW (FAISS)), we instead report the indexing time for rebuilding the index over the existing data. For Ada-ef, we adopt the incremental update strategy described in Section 6.3 to capture the effects of dataset changes. Specifically, for both insertion and deletion scenarios, we incrementally update (i) the dataset statistics and (ii) the ground truth of the sampled subset, and use these updated components to reconstruct the ef-estimation table. The corresponding update times are reported in Tables 4 and 6, respectively.

**Table 4: Insertion time (s).**

| Dataset | BS | HNSW (HNSWlib) | Ada-ef | Stats | Samp. | EF-Est. |
|---|---|---|---|---|---|---|
| DeepImage | 10% | 58.047 | 0.740 | 0.039 | 0.393 | 0.308 |
| DeepImage | 50% | 247.416 | 2.072 | 0.154 | 1.647 | 0.271 |
| MS MARCO V2.1 | 10% | 274.259 | 16.589 | 2.075 | 6.513 | 8.001 |
| MS MARCO V2.1 | 50% | 1213.185 | 46.175 | 9.613 | 28.318 | 8.244 |

**Table 5: Query time (QT) and recall of Ada-ef after insertion.**

| Dataset | Method | BS | QT (s) | Avg. Recall | P5 Recall | P1 Recall |
|---|---|---|---|---|---|---|
| DeepImage | Stale Ada-ef | 10% | 6.223 | 0.964 | 0.860 | 0.780 |
| DeepImage | Incr. Ada-ef | 10% | 7.463 | 0.973 | 0.900 | 0.830 |
| DeepImage | Stale Ada-ef | 50% | 4.828 | 0.945 | 0.820 | 0.720 |
| DeepImage | Incr. Ada-ef | 50% | 6.307 | 0.967 | 0.880 | 0.800 |
| DeepImage | Reco. Ada-ef | N/A | 6.437 | 0.968 | 0.880 | 0.810 |
| MS MARCO V2.1 | Stale Ada-ef | 10% | 15.262 | 0.952 | 0.865 | 0.789 |
| MS MARCO V2.1 | Incr. Ada-ef | 10% | 16.265 | 0.957 | 0.873 | 0.802 |
| MS MARCO V2.1 | Stale Ada-ef | 50% | 14.089 | 0.950 | 0.848 | 0.765 |
| MS MARCO V2.1 | Incr. Ada-ef | 50% | 14.230 | 0.951 | 0.833 | 0.625 |
| MS MARCO V2.1 | Reco. Ada-ef | N/A | 15.406 | 0.955 | 0.869 | 0.712 |

**Table 6: Deletion time of Ada-ef and indexing time (s).**

| Dataset | BS | HNSW (HNSWlib) | Ada-ef | Stats | Samp. | EF-Est. |
|---|---|---|---|---|---|---|
| DeepImage | 10% | 417.936 | 1.566 | 0.036 | 1.258 | 0.272 |
| DeepImage | 50% | 197.927 | 3.352 | 0.175 | 2.923 | 0.254 |
| MS MARCO V2.1 | 10% | 2079.652 | 13.538 | 2.209 | 4.270 | 7.059 |
| MS MARCO V2.1 | 50% | 1063.294 | 24.144 | 10.135 | 7.721 | 6.288 |

**Table 7: Query time (QT) and recall of Ada-ef after deletion.**

| Dataset | Method | BS | QT (s) | Avg. Recall | P5 Recall | P1 Recall |
|---|---|---|---|---|---|---|
| DeepImage | Stale Ada-ef | 10% | 6.297 | 0.969 | 0.890 | 0.810 |
| DeepImage | Incr. Ada-ef | 10% | 6.457 | 0.970 | 0.890 | 0.820 |
| DeepImage | Reco. Ada-ef | 10% | 6.089 | 0.966 | 0.870 | 0.790 |
| DeepImage | Stale Ada-ef | 50% | 6.025 | 0.979 | 0.920 | 0.860 |
| DeepImage | Incr. Ada-ef | 50% | 5.427 | 0.974 | 0.910 | 0.850 |
| DeepImage | Reco. Ada-ef | 50% | 4.373 | 0.960 | 0.860 | 0.780 |
| MS MARCO V2.1 | Stale Ada-ef | 10% | 16.235 | 0.954 | 0.870 | 0.708 |
| MS MARCO V2.1 | Incr. Ada-ef | 10% | 14.405 | 0.948 | 0.820 | 0.648 |
| MS MARCO V2.1 | Reco. Ada-ef | 10% | 15.723 | 0.953 | 0.864 | 0.782 |
| MS MARCO V2.1 | Stale Ada-ef | 50% | 16.821 | 0.960 | 0.886 | 0.678 |
| MS MARCO V2.1 | Incr. Ada-ef | 50% | 15.762 | 0.959 | 0.883 | 0.823 |
| MS MARCO V2.1 | Reco. Ada-ef | 50% | 14.793 | 0.956 | 0.871 | 0.770 |

To evaluate the search performance of Ada-ef, we consider three variants: (i) Stale Ada-ef, which uses the unmodified (pre-update) Ada-ef; (ii) Incr. Ada-ef, which applies the incrementally updated Ada-ef described above; and (iii) Reco. Ada-ef, which is fully recomputed from scratch. Their search performance results are reported in Tables 5 and 7, respectively.

As shown in Tables 4 and 6, the update time required by Ada-ef is modest compared to the HNSW index update time in the insertion scenario and the index rebuilding time in the deletion scenario.

For search performance, in the insertion case, Stale Ada-ef maintains robust performance but may fall slightly below the target recall, particularly with large batch sizes (*e.g.*, 50% insertion on Deep-Image). In the deletion case, Stale Ada-ef consistently achieves recall above the target of 0.95 across all datasets, though at the expense of slightly higher search time. The incrementally updated version, Incr. Ada-ef, effectively resolves this issue, maintaining target recall while avoiding over-searching, and achieves performance comparable to the fully recomputed Reco. Ada-ef. In summary,

**Table 8: Impact of the distance list size |D| on the offline computation time and online query time of Ada-ef.**

| Dataset | \|D\| | EF-Est. (s) | QT (s) | Avg. Recall | P5 Recall | P1 Recall |
|---|---|---|---|---|---|---|
| DeepImage | 1-hop | 0.349 | 6.524 | 0.968 | 0.870 | 0.770 |
| DeepImage | 2-hop | 0.282 | 6.902 | 0.968 | 0.880 | 0.810 |
| DeepImage | 3-hop | 3.741 | 92.410 | 0.998 | 0.990 | 0.960 |
| MS MARCO V2.1 | 1-hop | 3.081 | 12.890 | 0.943 | 0.783 | 0.626 |
| MS MARCO V2.1 | 2-hop | 3.171 | 15.622 | 0.955 | 0.869 | 0.712 |
| MS MARCO V2.1 | 3-hop | 7.540 | 30.871 | 0.968 | 0.849 | 0.702 |

**Table 9: Impact of the number of sampled data vectors on the offline computation time and online query time of Ada-ef.**

| Dataset | Num. | Samp. (s) | EF-Est. (s) | QT (s) | Avg. Recall | P5 Recall | P1 Recall |
|---|---|---|---|---|---|---|---|
| DeepImage | 200 | 6.827 | 0.273 | 6.420 | 0.968 | 0.880 | 0.810 |
| DeepImage | 1000 | 23.704 | 1.463 | 6.342 | 0.966 | 0.880 | 0.820 |
| DeepImage | 5000 | 99.100 | 8.657 | 4.713 | 0.940 | 0.770 | 0.650 |
| MS MARCO V2.1 | 200 | 84.240 | 3.155 | 15.348 | 0.955 | 0.869 | 0.712 |
| MS MARCO V2.1 | 1000 | 102.390 | 14.455 | 13.401 | 0.948 | 0.824 | 0.636 |
| MS MARCO V2.1 | 5000 | 260.469 | 83.036 | 13.610 | 0.950 | 0.845 | 0.662 |

**Table 10: Impact of the design of weighted function on the offline computation time and online query time of Ada-ef.**

| Dataset | Decay | EF-Est. (s) | QT (s) | Avg. Recall | P5 Recall | P1 Recall |
|---|---|---|---|---|---|---|
| DeepImage | None | 0.649 | 5.707 | 0.951 | 0.800 | 0.680 |
| DeepImage | Linear | 0.301 | 5.419 | 0.946 | 0.790 | 0.660 |
| DeepImage | Exp | 0.286 | 6.950 | 0.968 | 0.880 | 0.81 |
| MS MARCO V2.1 | None | 2.663 | 9.594 | 0.919 | 0.732 | 0.563 |
| MS MARCO V2.1 | Linear | 2.723 | 12.858 | 0.944 | 0.804 | 0.659 |
| MS MARCO V2.1 | Exp | 3.239 | 15.987 | 0.955 | 0.869 | 0.712 |

Stale Ada-ef performs robustly under small batch updates, whereas Incr. Ada-ef is preferable for large batch updates, offering a better balance between recall stability and computational efficiency.

## 7.6 Ablation Study on Component Designs

We conduct ablation studies on two representative datasets, Deep-Image and MS MARCO V2.1, to analyze the effectiveness of different component designs in Ada-ef.

The first ablation study evaluates the impact of the distance list size |D|, considering three configurations: neighbors within 1-hop, 2-hop, and 3-hop. The results are shown in Table 8. In general, the 3-hop configuration substantially increases both query latency and EF-Est. table computation time, as the search must traverse all neighbors within three hops from the entry point on the base layer. The 1-hop configuration suffers from low recall, particularly on high-dimensional datasets such as MS MARCO V2.1. The 2-hop design achieves the best balance between recall and latency, and is thus adopted as the default configuration.

The second study investigates the effect of sampling size, with three settings: 200, 1000, and 5000 sampled vectors. The results are presented in Table 9. Larger sampling sizes significantly increase offline computation time without consistently improving online search performance. Sampling 200 vectors provides the best overall trade-off between accuracy and efficiency.

The third study examines the influence of the decay function, which assigns weights used in Eq. (6). We evaluate three alternatives: no decay function (None), a linear decay function (Linear), and an exponential decay function (Exp). The results are shown in Table 10. The None setting assigns equal weights to all bins in

Eq. (6), failing to distinguish between easy and hard queries and resulting in lower recall, especially on MS MARCO V2.1. The Linear function gives higher weights to closer neighbors, partially improving discrimination but still less effectively than Exp. The Exponential decay function yields the best performance, avoiding both over- and under-searching, as evidenced by improvements in average recall as well as 1st- and 5th-percentile recall.

## 8 Conclusion

HNSW has been widely adopted in practice for efficient ANNS at scale. However, its search performance is highly sensitive to the choice of the ef parameter, which governs the trade-off between efficiency and recall. The commonly used fixed ef configuration suffers from practical limitations, including lack of recall guarantees and issues related to over- or under-searching across diverse queries. In this paper, we identify these limitations and propose Ada-ef, a distribution-aware approach that adaptively assigns an ef value to each incoming query at runtime. Ada-ef enables the search to approximately meet user-specified target recall while minimizing unnecessary computation, thus mitigating both under- and over-searching. Ada-ef is seamlessly integrable with existing HNSW indexes, requiring no modification to the index structure, and is friendly to updates, making it suitable for real-time and evolving workloads. At the core of Ada-ef lies a theoretical foundation based on similarity distance distribution, along with a tailored query scoring model designed to effectively characterize query difficulty. On real-world embeddings produced by advanced Transformer models from OpenAI and Cohere, Ada-ef attains the declarative recall while mitigating over- and under-searching. Compared to state-of-the-art learning-based adaptive baselines, Ada-ef achieves up to 4× reduction in online query latency, 50× reduction in offline computation, and 100× reduction in offline memory usage.

## References

[1] [n.d.]. *FAISS: A Library for Efficient Similarity Search and Clustering of Dense Vectors.* https://github.com/facebookresearch/faiss Accessed: 2025-10-21.
[2] [n.d.]. *Hnswlib: Fast Approximate Nearest Neighbor Search.* https://github.com/nmslib/hnswlib Accessed: 2025-10-21.
[3] 2006. *Eigen is a C++ template library for linear algebra: matrices, vectors, numerical solvers, and related algorithms.* https://eigen.tuxfamily.org/index.php?title=Main_Page Accessed: 2025-01-10.
[4] 2019. *ann-benchmarks.* https://github.com/erikbern/ann-benchmarks
[5] 2019. *Pinecone.* https://www.pinecone.io/ Accessed: 2025-05-19.
[6] 2019. *Weaviate.* https://weaviate.io/ Accessed: 2025-03-19.
[7] 2020. *MS MARCO.* https://microsoft.github.io/msmarco
[8] 2021. *OpenSearch Project.* https://github.com/opensearch-project
[9] 2021. *pgvector: Open-source vector similarity search for Postgres.*
[10] 2022. *HnswGraph (Lucene 9.8.0 core API).* https://lucene.apache.org/core/9_8_0/core/org/apache/lucene/util/hnsw/HnswGraph.html
[11] 2024. *Anserini: OpenAI-ada2 Embeddings for MS MARCO Passage.* https://github.com/castorini/anserini/blob/master/docs/experiments-msmarco-passage-openai-ada2.md?utm_source=chatgpt.com Accessed: 2025-03-20.
[12] 2024. *SISAP 2024 Indexing Challenge.* https://sisap-challenges.github.io/2025/index.html
[13] 2024. *TREC 2024 RAG Corpus.* https://trec-rag.github.io/annoucements/2024-corpus-finalization
[14] 2024. *TREC-RAG 2024 Corpus (MSMARCO 2.1) - Encoded with Cohere Embed English v3.* https://huggingface.co/datasets/Cohere/msmarco-v2.1-embed-english-v3 Accessed: 2024-12-10.
[15] Thomas D. Ahle, Martin Aumüller, and Rasmus Pagh. 2017. Parameter-free locality sensitive hashing for spherical range reporting. In *Proc. 28th Annual ACM-SIAM Symp. on Discrete Algorithms (SODA '17)*. 239–256.
[16] Laurent Amsaleg, Oussama Chelly, Teddy Furon, Stéphane Girard, Michael E. Houle, Ken-ichi Kawarabayashi, and Michael Nett. 2015. Estimating Local

[17] Intrinsic Dimensionality. In *Proc. 21st ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*. 29–38. https://doi.org/10.1145/2783258.2783405
[18] Fabien André, Anne-Marie Kermarrec, and Nicolas Le Scouarnec. 2015. Cache locality is not enough: high-performance nearest neighbor search with product quantization fast scan. *Proc. VLDB Endowment* 9, 4 (Dec. 2015), 288–299. https://doi.org/10.14778/2856318.2856324
[18] Akhil Arora, Sakshi Sinha, Piyush Kumar, and Arnab Bhattacharya. 2018. HD-index: pushing the scalability-accuracy boundary for approximate kNN search in high-dimensional spaces. *Proc. VLDB Endowment* 11, 8 (April 2018), 906–919. https://doi.org/10.14778/3204028.3204034
[19] Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. 2016. A Latent Variable Model Approach to PMI-based Word Embeddings. *Transactions of the Association for Computational Linguistics* 4 (2016), 385–399. https://doi.org/10.1162/tacl_a_00106
[20] Martin Aumüller, Erik Bernhardsson, and Alexander Faithfull. 2018. ANN-Benchmarks: A Benchmarking Tool for Approximate Nearest Neighbor Algorithms. arXiv:1807.05614 [cs.IR] https://arxiv.org/abs/1807.05614
[21] Ilias Azizi, Karima Echihabi, and Themis Palpanas. 2023. ELPIS: Graph-Based Similarity Search for Scalable Data Science. *Proc. VLDB Endowment* 16, 6 (Feb. 2023), 1548–1559. https://doi.org/10.14778/3583140.3583166
[22] Ilias Azizi, Karima Echihabi, and Themis Palpanas. 2025. Graph-Based Vector Search: An Experimental Evaluation of the State-of-the-Art. *Proc. ACM Manag. Data* 3, 1, Article 43 (Feb. 2025), 31 pages. https://doi.org/10.1145/3709693
[23] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450* (2016).
[24] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. 2018. MS MARCO: A Human Generated MAchine Reading COmprehension Dataset. arXiv:1611.09268 [cs.CL] https://arxiv.org/abs/1611.09268
[25] Erik Bernhardsson. 2017. *Annoy: approximate nearest neighbors in C++/Python optimized for memory usage and loading/saving to disk.* https://github.com/spotify/annoy
[26] Simion-Vlad Bogolin, Ioana Croitoru, Hailin Jin, Yang Liu, and Samuel Albanie. 2022. Cross Modal Retrieval with Querybank Normalisation. In *Proc. IEEE Int. Conf. on Computer Vision and Pattern Recognition*. 5184–5195. https://doi.org/10.1109/CVPR52688.2022.00513
[27] Manos Chatzakis, Yannis Papakonstantinou, and Themis Palpanas. 2025. DARTH: Declarative Recall Through Early Termination for Approximate Nearest Neighbor Search. *Proc. ACM Manag. Data* 3, 4 (Sept. 2025), 26. https://doi.org/10.1145/3749160
[28] Cheng Chen, Chenzhe Jin, Yunan Zhang, Sasha Podolsky, Chun Wu, Szu-Po Wang, Eric Hanson, Zhou Sun, Robert Walzer, and Jianguo Wang. 2024. SingleStore-V: An Integrated Vector Database System in SingleStore. *Proc. VLDB Endowment* 17, 12 (Aug. 2024), 3772–3785. https://doi.org/10.14778/3685800.3685805
[29] Patrick Chen, Wei-Cheng Chang, Jyun-Yu Jiang, Hsiang-Fu Yu, Inderjit Dhillon, and Cho-Jui Hsieh. 2023. FINGER: Fast Inference for Graph-based Approximate Nearest Neighbor Search. In *ProC. Web Conf. 2023*. 3225–3235. https://doi.org/10.1145/3543507.3583318
[30] Qi Chen, Haidong Wang, Mingqin Li, Gang Ren, Scarlett Li, Jeffery Zhu, Jason Li, Chuanjie Liu, Lintao Zhang, and Jingdong Wang. 2018. SPTAG: A library for fast approximate nearest neighbor search. https://github.com/Microsoft/SPTAG
[31] 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 Neighbor Search. In *Advances in Neural Information Processing Syst., Proc. 35th Annual Conf. on Neural Information Processing Syst.*
[32] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 4171–4186. https://doi.org/10.18653/v1/N19-1423
[33] Wei Dong. 2011. *KGraph: A Library for Approximate Nearest Neighbor Search.* https://github.com/aaalgo/kgraph
[34] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *Proc. 9th Int. Conf. on Learning Representations*. https://openreview.net/forum?id=YicbFdNTTy
[35] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The Faiss library. (2024). arXiv:2401.08281 [cs.LG]
[36] Owen Elliott. 2024. *Understanding Recall in HNSW Search.* https://www.marqo.ai/blog/understanding-recall-in-hnsw-search
[37] Kawin Ethayarajh. 2019. How Contextual are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings. 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).* 55–65. https://doi.org/10.18653/v1/D19-1006

[38] Xiyang Feng, Guodong Jin, Ziyi Chen, Chang Liu, and Semih Salihoğlu. 2023. Kùzu graph database management system. In *The Conf. on Innovative Data Systems Research*, Vol. 7. 25–35.

[39] Jerome H. Friedman. 2001. Greedy Function Approximation: A Gradient Boosting Machine. *The Annals of Statistics* 29, 5 (2001), 1189–1232. http://www.jstor.org/stable/2699986

[40] Cong Fu and Deng Cai. 2016. EFANNA : An Extremely Fast Approximate Nearest Neighbor Search Algorithm Based on kNN Graph. arXiv:1609.07228 [cs.CV] https://arxiv.org/abs/1609.07228

[41] Cong Fu, Changxu Wang, and Deng Cai. 2022. High Dimensional Similarity Search With Satellite System Graph: Efficiency, Scalability, and Unindexed Query Compatibility. *IEEE Trans. Pattern Analy. Machine Intell.* 44, 8 (2022), 4139–4150. https://doi.org/10.1109/TPAMI.2021.3067706

[42] Cong Fu, Chao Xiang, Changxu Wang, and Deng Cai. 2019. Fast approximate nearest neighbor search with the navigating spreading-out graph. *Proc. VLDB Endowment* 12, 5 (Jan. 2019), 461–474. https://doi.org/10.14778/3303753.3303754

[43] Max Gabrielsson. 2024. *Vector Similarity Search in DuckDB.* https://duckdb.org/2024/05/03/vector-similarity-search-vss.html

[44] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997* 2 (2023).

[45] Long Gong, Huayi Wang, Mitsunori Ogihara, and Jun Xu. 2020. iDEC: indexable distance estimating codes for approximate nearest neighbor search. *Proc. VLDB Endowment* 13, 9 (May 2020), 1483–1497. https://doi.org/10.14778/3397230.3397243

[46] Ryan Greene, Ted Sanders, Lilian Weng, and Arvind Neelakantan. 2022. *New and improved embedding model: text-embedding-ada-002.* https://openai.com/index/new-and-improved-embedding-model

[47] Ben Harwood and Tom Drummond. 2016. FANNG: Fast Approximate Nearest Neighbour Graphs. In *Proc. IEEE Int. Conf. on Computer Vision and Pattern Recognition.*

[48] Junfeng He, Sanjiv Kumar, and Shih-Fu Chang. 2012. On the difficulty of nearest neighbor search. In *Proc. 29th Int. Conf. on Machine Learning.* 41–48.

[49] Qiang Huang, Jianlin Feng, Yikai Zhang, Qiong Fang, and Wilfred Ng. 2015. Query-aware locality-sensitive hashing for approximate nearest neighbor search. *Proc. VLDB Endowment* 9, 1 (Sept. 2015), 1–12. https://doi.org/10.14778/2850469.2850470

[50] 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. https://doi.org/10.1145/276698.276876

[51] Masajiro Iwasaki. 2015. *Neighborhood Graph and Tree for Indexing Highdimensional Data.* Yahoo Japan Corporation. https://github.com/yahoojapan/NGT.

[52] Omid Jafari, Preeti Maurya, Parth Nagarkar, Khandker Mushfiqul Islam, and Chidambaram Crushev. 2021. A Survey on Locality Sensitive Hashing Algorithms and their Applications. arXiv:2102.08942 [cs.DB] https://arxiv.org/abs/2102.08942

[53] Suhas Jayaram Subramanya, Fnu Devvrit, Harsha Vardhan Simhadri, Ravishankar Krishnawamy, and Rohan Kadekodi. 2019. DiskANN: Fast Accurate Billion-point Nearest Neighbor Search on a Single Node. In *Advances in Neural Information Processing Syst., Proc. 33rd Annual Conf. on Neural Information Processing Syst.*, Vol. 32. https://proceedings.neurips.cc/paper_files/paper/2019/file/09853c7fb1d3f8ee67a61b6bf4a7f8e6-Paper.pdf

[54] Zhongming Jin, Debing Zhang, Yao Hu, Shiding Lin, Deng Cai, and Xiaofei He. 2014. Fast and Accurate Hashing Via Iterative Nearest Neighbors Expansion. *IEEE Trans. Cybernetics* 44, 11 (2014), 2167–2177. https://doi.org/10.1109/TCYB.2014.2302018

[55] Herve Jégou, Matthijs Douze, and Cordelia Schmid. 2011. Product Quantization for Nearest Neighbor Search. *IEEE Trans. Pattern Analy. Machine Intell.* 33, 1 (2011), 117–128. https://doi.org/10.1109/TPAMI.2010.57

[56] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP).* 6769–6781. https://doi.org/10.18653/v1/2020.emnlp-main.550

[57] Achim Klenke. 2020. *Characteristic Functions and the Central Limit Theorem.* Springer International Publishing, Cham, 327–366. https://doi.org/10.1007/978-3-030-56402-5_15

[58] Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020. On the Sentence Embeddings from Pre-trained Language Models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP).* 9119–9130. https://doi.org/10.18653/v1/2020.emnlp-main.733

[59] Conglong Li, Minjia Zhang, David G. Andersen, and Yuxiong He. 2020. Improving Approximate Nearest Neighbor Search through Learned Adaptive Early

[60] Haitao Li, Qingyao Ai, Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Zheng Liu, and Zhao Cao. 2023. Constructing Tree-based Index for Efficient and Effective Dense Retrieval. In *Proc. 46th Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval.* 131–140. https://doi.org/10.1145/3539618.3591651

[61] Wen Li, Ying Zhang, Yifang Sun, Wei Wang, Mingjie Li, Wenjie Zhang, and Xuemin Lin. 2020. Approximate Nearest Neighbor Search on High Dimensional Data — Experiments, Analyses, and Improvement. *IEEE Trans. Knowl. and Data Eng.* 32, 8 (2020), 1475–1488. https://doi.org/10.1109/TKDE.2019.2909204

[62] Shige Liu, Zhifang Zeng, Li Chen, Adil Ainihaer, Arun Ramasami, Songting Chen, Yu Xu, Mingxi Wu, and Jianguo Wang. 2025. TigerVector: Supporting Vector Search in Graph Databases for Advanced RAGs. arXiv:2501.11216 [cs.DB] https://arxiv.org/abs/2501.11216

[63] Kejing Lu, Hongya Wang, Wei Wang, and Mineichi Kudo. 2020. VHP: approximate nearest neighbor search via virtual hypersphere partitioning. *Proc. VLDB Endowment* 13, 9 (May 2020), 1443–1455. https://doi.org/10.14778/3397230.3397240

[64] 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. https://doi.org/10.1016/j.is.2013.10.006

[65] Yu A. Malkov and D. A. Yashunin. 2020. Efficient and Robust Approximate Nearest Neighbor Search Using Hierarchical Navigable Small World Graphs. *IEEE Trans. Pattern Analy. Machine Intell.* 42, 4 (2020), 824–836. https://doi.org/10.1109/TPAMI.2018.2889473

[66] Xiangrui Meng, Joseph Bradley, Burak Yavuz, Evan Sparks, Shivaram Venkataraman, Davies Liu, Jeremy Freeman, DB Tsai, Manish Amde, Sean Owen, Doris Xin, Reynold Xin, Michael J. Franklin, Reza Zadeh, Matei Zaharia, and Ameet Talwalkar. 2016. MLlib: Machine Learning in Apache Spark. *Journal of Machine Learning Research* 17, 34 (2016), 1–7. http://jmlr.org/papers/v17/15-237.html

[67] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Syst., Proc. 27th Annual Conf. on Neural Information Processing Syst.*, C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger (Eds.), Vol. 26. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf

[68] Jiaqi Mu and Pramod Viswanath. 2018. All-but-the-Top: Simple and Effective Postprocessing for Word Representations. In *Proc. 6th Int. Conf. on Learning Representations.*

[69] Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav Shyam, Boris Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David Schnurr, Felipe Petroski Such, Kenny Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov, Joanne Jang, Peter Welinder, and Lilian Weng. 2022. Text and Code Embeddings by Contrastive Pre-Training. arXiv:2201.10005 [cs.CL] https://arxiv.org/abs/2201.10005

[70] Beatrix Miranda Ginn Nielsen and Lars Kai Hansen. 2024. Hubness Reduction Improves Sentence-BERT Semantic Spaces. In *Proceedings of the 5th Northern Lights Deep Learning Conference (NLDL) (Proceedings of Machine Learning Research)*, Tetiana Lutchyn, Adín Ramírez Rivera, and Benjamin Ricaud (Eds.), Vol. 233. 181–204.

[71] Hiroyuki Ootomo, Akira Naruse, Corey Nolet, Ray Wang, Tamas Feher, and Yong Wang. 2024. CAGRA: Highly Parallel Graph Construction and Approximate Nearest Neighbor Search for GPUs. In *Proc. 36th Int. Conf. on Data Engineering.* 4236–4247. https://doi.org/10.1109/ICDE60146.2024.00323

[72] James Jie Pan, Jianguo Wang, and Guoliang Li. 2024. Survey of vector database management systems. *VLDB J.* 33, 5 (July 2024), 1591–1615. https://doi.org/10.1007/s00778-024-00864-x

[73] Zhibin Pan, Liangzhuang Wang, Yang Wang, and Yuchen Liu. 2020. Product quantization with dual codebooks for approximate nearest neighbor search. *Neurocomputing* 401 (2020), 59–68. https://doi.org/10.1016/j.neucom.2020.03.016

[74] Zhen Peng, Minjia Zhang, Kai Li, Ruoming Jin, and Bin Ren. 2023. iQAN: Fast and Accurate Vector Search with Efficient Intra-Query Parallelism on Multi-Core Architectures. In *Proc. 28th ACM SIGPLAN Symp. on Principles and Practice of Parallel Programming.* 313–328. https://doi.org/10.1145/3572848.3577527

[75] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global Vectors for Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Alessandro Moschitti, Bo Pang, and Walter Daelemans (Eds.). Association for Computational Linguistics, Doha, Qatar, 1532–1543. https://doi.org/10.3115/v1/D14-1162

[76] Mark Raasveldt and Hannes Mühleisen. 2019. DuckDB: an Embeddable Analytical Database. In *Proc. ACM SIGMOD Int. Conf. on Management of Data.* 1981–1984. https://doi.org/10.1145/3299869.3320212

[77] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark,

Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proc. 38th Int. Conf. on Machine Learning.* https://api.semanticscholar.org/CorpusID:231591445

[78] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proc. 38th Int. Conf. on Machine Learning,* Vol. 139. 8748–8763. https://proceedings.mlr.press/v139/radford21a.html

[79] Nils Reimers and Iryna Gurevych. 2019. 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).* 3982–3992. https://doi.org/10.18653/v1/D19-1410

[80] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. LAION-400M: Open Dataset of CLIP-Filtered 400 Million Image-Text Pairs. arXiv:2111.02114 [cs.CV] https://arxiv.org/abs/2111.02114

[81] Chanop Silpa-Anan and Richard Hartley. 2008. Optimised KD-trees for fast image descriptor matching. In *Proc. IEEE Int. Conf. on Computer Vision and Pattern Recognition.* 1–8. https://doi.org/10.1109/CVPR.2008.4587638

[82] Tommaso Teofili and Jimmy Lin. 2025. Patience in Proximity: A Simple Early Termination Strategy for HNSW Graph Traversal Approximate k-Nearest Neighbor Search. In *Proc. 47th European Conf. on IR Research.* 401–407. https://doi.org/10.1007/978-3-031-88714-7_39

[83] Bing Tian, Haikun Liu, Yuhang Tang, Shihai Xiao, Zhuohui Duan, Xiaofei Liao, Hai Jin, Xuecang Zhang, Junhua Zhu, and Yu Zhang. 2025. Towards High-throughput and Low-latency Billion-scale Vector Search via CPU/GPU Collaborative Filtering and Re-ranking. In *Proc. 23rd USENIX Conf. on File and Storage Technologies.* 171–185. https://www.usenix.org/conference/fast25/presentation/tian-bing

[84] Julie Tibshirani. 2022. *Introducing approximate nearest neighbor search in Elasticsearch 8.0.* https://www.elastic.co/blog/introducing-approximate-nearest-neighbor-search-in-elasticsearch-8-0

[85] Kirill Tyshchuk, Polina Karpikova, Andrew Spiridonov, Anastasiia Prutianova, Anton Razzhigaev, and Alexander Panchenko. 2023. On Isotropy of Multimodal Embeddings. *Inf.* 14, 7 (2023), 392.

[86] Javier Vargas Muñoz, Marcos A. Gonçalves, Zanoni Dias, and Ricardo da S. Torres. 2019. Hierarchical Clustering-Based Graphs for Large Scale Approximate Nearest Neighbor Search. *Pattern Recognition* 96 (2019), 106970. https://doi.org/10.1016/j.patcog.2019.106970

[87] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Syst., Proc. 31st Annual Conf. on Neural Information Processing Syst.,* I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

[88] Jianguo Wang, Xiaomeng Yi, Rentong Guo, Hai Jin, Peng Xu, Shengjun Li, Xiangyu Wang, Xiangzhou Guo, Chengming Li, Xiaohai Xu, Kun Yu, Yuxing Yuan, Yinghao Zou, Jiquan Long, Yudong Cai, Zhenxiang Li, Zhifeng Zhang, Yihua Mo, Jun Gu, Ruiyi Jiang, Yi Wei, and Charles Xie. 2021. Milvus: A Purpose-Built Vector Data Management System. In *Proc. ACM SIGMOD Int. Conf. on Management of Data.* 2614–2627. https://doi.org/10.1145/3448016.3457550

[89] Mengzhao Wang, Haotian Wu, Xiangyu Ke, Yunjun Gao, Yifan Zhu, and Wenchao Zhou. 2025. Accelerating Graph Indexing for ANNS on Modern CPUs. *arXiv preprint arXiv:2502.18113* (2025).

[90] Mengzhao Wang, Weizhi Xu, Xiaomeng Yi, Songlin Wu, Zhangyang Peng, Xiangyu Ke, Yunjun Gao, Xiaoliang Xu, Rentong Guo, and Charles Xie. 2024. Starling: An I/O-Efficient Disk-Resident Graph Index Framework for High-Dimensional Vector Similarity Search on Data Segment. *Proc. ACM Manag. Data* 2, 1, Article 14 (March 2024), 27 pages. https://doi.org/10.1145/3639269

[91] Mengzhao Wang, Xiaoliang Xu, Qiang Yue, and Yuxiang Wang. 2021. A comprehensive survey and experimental comparison of graph-based approximate nearest neighbor search. *Proc. VLDB Endowment* 14, 11 (July 2021), 1964–1978. https://doi.org/10.14778/3476249.3476255

[92] Zeyu Wang, Peng Wang, Themis Palpanas, and Wei Wang. 2023. Graph- and Tree-based Indexes for High-dimensional Vector Similarity Search: Analyses, Comparisons, and Future Directions. *IEEE Data Eng. Bull.* 46 (2023), 3–21. https://api.semanticscholar.org/CorpusID:265509541

[93] Zeyu Wang, Qitong Wang, Xiaoxing Cheng, Peng Wang, Themis Palpanas, and Wei Wang. 2024. Steiner-Hardness: A Query Hardness Measure for Graph-Based ANN Indexes. *Proc. VLDB Endowment* 17, 13 (2024), 4668–4682. https://doi.org/10.14778/3704965.3704974

[94] Artem Babenko Yandex and Victor Lempitsky. 2016. Efficient Indexing of Billion-Scale Datasets of Deep Descriptors. In *Proc. IEEE Int. Conf. on Computer Vision and Pattern Recognition.* 2055–2063. https://doi.org/10.1109/CVPR.2016.226

[95] Kaixiang Yang, Hongya Wang, Bo Xu, Wei Wang, Yingyuan Xiao, Ming Du, and Junfeng Zhou. 2021. Tao: A Learning Framework for Adaptive Nearest Neighbor Search using Static Features Only. arXiv:2110.00696 [cs.DB] https://arxiv.org/abs/2110.00696

[96] Peilin Yang, Hui Fang, and Jimmy Lin. 2017. Anserini: Enabling the Use of Lucene for Information Retrieval Research. In *Proc. 40th Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval.* 1253–1256. https://doi.org/10.1145/3077136.3080721

[97] Wen Yang, Tao Li, Gai Fang, and Hong Wei. 2020. PASE: PostgreSQL Ultra-High-Dimensional Approximate Nearest Neighbor Search Extension. In *Proc. ACM SIGMOD Int. Conf. on Management of Data.* 2241–2253. https://doi.org/10.1145/3318464.3386131

[98] Sho Yokoi, Han Bao, Hiroto Kurita, and Hidetoshi Shimodaira. 2024. Zipfian Whitening. In *Advances in Neural Information Processing Syst., Proc. 38th Annual Conf. on Neural Information Processing Syst.*

[99] Kisung You. 2025. Semantics at an Angle: When Cosine Similarity Works Until It Doesn't. arXiv:2504.16318 [cs.LG] https://arxiv.org/abs/2504.16318

[100] Qianxi Zhang, Shuotao Xu, Qi Chen, Guoxin Sui, Jiadong Xie, Zhizhen Cai, Yaoqi Chen, Yinxuan He, Yuqing Yang, Fan Yang, et al. 2023. {VBASE}: Unifying Online Vector Similarity Search and Relational Queries via Relaxed Monotonicity. In *Proc. 17th USENIX Symp. on Operating System Design and Implementation.*

[101] Xi Zhao, Bolong Zheng, Xiaomeng Yi, Xiaofan Luan, Charles Xie, Xiaofang Zhou, and Christian S. Jensen. 2023. FARGO: Fast Maximum Inner Product Search via Global Multi-Probing. *Proc. VLDB Endowment* 16, 5 (Jan. 2023), 1100–1112. https://doi.org/10.14778/3579075.3579084

## A Pseudocode of Adaptive EF Search

---

**Algorithm 2** ADAPTIVE-EF-SEARCH($\mathbf{q}, ep, r$)

---

**Input:** query vector $\mathbf{q}$, enter point $ep$, expected recall $r$
**Output:** closest neighbors to $\mathbf{q}$

1: $S \leftarrow ep$                 ▷ set of visited elements
2: $C \leftarrow ep$                 ▷ set of candidates
3: $W \leftarrow ep$      ▷ dynamic list of found nearest neighbors
4: $\mathsf{ef} \leftarrow \infty$        ▷ no limitation of ef at the beginning
5: $D \leftarrow$ get distance between $ep$ and $\mathbf{q}$   ▷ sampled distance list
6: **while** $|C| > 0$ **do**
7:     $\mathbf{c} \leftarrow$ extract from $C$ the nearest element to $\mathbf{q}$
8:     $\mathbf{f} \leftarrow$ get from $W$ the furthest element to $\mathbf{q}$
9:     **if** $dist(\mathbf{c}, \mathbf{q}) > dist(\mathbf{f}, \mathbf{q})$ **then**
10:         **break**
11:     **for** each $\mathbf{e} \in$ neighborhood($\mathbf{c}$) **do**     ▷ update $C$ and $W$
12:         **if** $\mathbf{e} \notin S$ **then**
13:             $S \leftarrow S \cup \mathbf{e}$
14:             $\mathbf{f} \leftarrow$ get from $W$ the furthest element to $\mathbf{q}$
15:             **if** $dist(\mathbf{e}, \mathbf{q}) < dist(\mathbf{f}, \mathbf{q})$ or $|W| < \mathsf{ef}$ **then**
16:                 $C \leftarrow C \cup \mathbf{e}$
17:                 $W \leftarrow W \cup \mathbf{e}$
18:                 **if** $|W| > ef$ **then**
19:                     remove from $W$ the furthest element to $\mathbf{q}$
20:         **if** $|D| < l$ **then**         ▷ collect distances
21:             $D \leftarrow dist(\mathbf{e}, \mathbf{q})$
22:         **if** $|D| = l$ **then**       ▷ reach distance limit
23:             $\mathsf{ef} \leftarrow$ ESTIMATE-EF($\mathbf{q}, D, r$)
24:             **while** $|W| > \mathsf{ef}$ **do**
25:                 remove from $W$ the furthest element
26: **return** $W$

---

Algorithm 2 presents the adaptive-ef search algorithm, which extends the standard HNSW search described in Section 3, with several key distinctions elaborated below. First, the input to Algorithm 2 does not require a predefined ef value; instead, ef is initially set to $\infty$, allowing the traversal to visit any nodes around

the entry point $ep$. Additionally, the algorithm initializes a sampled distance list $D$ to record the distances between visited nodes and the query vector $\mathbf{q}$. The size of $D$ is bounded by $l$ that corresponds to the number of nodes reachable within 2-hop from the entry point $ep$ at the base layer. The primary distinction from the original HNSW search lies between lines 20 and 25 of Algorithm 2. Once $l$ distances have been collected, the algorithm invokes the ESTIMATE-EF function, which takes the query vector $\mathbf{q}$, the list $D$, and the user-specified recall $r$ as input, and returns a dynamically estimated ef value. If the estimated ef exceeds $l$, the dynamic list $W$ which maintains the current nearest neighbors will be truncated to match the estimated ef.

## B  Relaxing Identical Distribution

The classical CLT can be extended to sequences of non-identically distributed random variables under certain conditions—most notably, Lindeberg's condition [57]. Consider the inner product between a query vector and a data vector, $\mathbf{q} \cdot \mathbf{v} = \sum_{i=1}^{d} q_i v_i$. For convenience, let each term $q_i v_i$ be represented by a random variable $X_i$. Suppose each $X_i$ has mean $\mu_i$ and variance $\sigma_i^2$. Then, Lindeberg's condition for the sequence $(X_1, \ldots, X_d)$ is satisfied if, for every $\varepsilon > 0$,

$$\lim_{d \to \infty} \frac{1}{s_d^2} \sum_{i=1}^{d} \mathrm{E}\left[ (X_i - \mu_i)^2 \cdot \mathbf{1}_{\{|X_i - \mu_i| > \varepsilon s_d\}} \right] = 0,$$

where $s_d^2 = \sum_{i=1}^{d} \sigma_i^2$, and $\mathbf{1}_{\{|X_i - \mu_i| > \varepsilon s_d\}}$ is the indicator function, equal to 1 when the condition inside holds, and 0 otherwise. Lindeberg's condition ensures that *no single $X_i$ contributes a disproportionately large amount to the sum $\sum_{i=1}^{d} X_i$*. In practice, modern embedding models—especially those based on Transformers—often apply layer normalization [23, 87] during training that regulates the scale of embedding components, resulting in components with approximately balanced magnitudes. This empirical property supports the assumption that real-world embedding vectors are likely to satisfy Lindeberg's condition.

## C  Example of Query Score

Consider the case $m = 5$ and $\delta = 0.001$, and using cosine distance as the similarity metric. Assume that for a given query vector $\mathbf{q}$, the estimated FDL distribution is parameterized by $\mu = 0.936$ and $\sigma = 0.0739$, as obtained from Equation 3. Under this setting, the first bin $b_1$ contains all collected distances in $D$ that are smaller than $0.936 + 0.0739 \times \Phi^{-1}(0.001) = 0.7076$, while bin $b_2$ includes distances between 0.7076 and $0.7233 = 0.936 + 0.0739 \times \Phi^{-1}(0.002)$. The remaining three bins are defined similarly based on the corresponding quantile thresholds. Suppose the total number of distances is $|D| = 100$, and the bin counts are $c_1 = 90$, $c_2 = 5$, and $c_3 = 5$, with $c_4 = c_5 = 0$. The query score for $\mathbf{q}$ is then computed as: $\mathrm{Score}(\mathbf{q}) = 0.9 \times 100.0 + 0.05 \times 36.79 + 0.05 \times 13.53 = 92.516$. This high score suggests that $\mathbf{q}$ could be an "easy" query, requiring a small ef to achieve target recall. Indeed, 90% of the first 100 visited nodes have distances smaller than the 0.001 quantile of the FDL distribution, indicating a high concentration of strong candidates early in the search.