# Down with the Hierarchy: The 'H' in HNSW Stands for "Hubs"

Blaise Munyampirwa
Independent Researcher
Mountain View, CA
blaisemunyampirwa@gmail.com

Vihan Lakshman
MIT CSAIL
Cambridge, MA
vihan@mit.edu

Benjamin Coleman
Google DeepMind
Mountain View, CA
colemanben@google.com

## ABSTRACT

Driven by recent breakthrough advances in neural representation learning, approximate near-neighbor (ANN) search over vector embeddings has emerged as a critical computational workload. With the introduction of the seminal Hierarchical Navigable Small World (HNSW) algorithm, graph-based indexes have established themselves as the overwhelmingly dominant paradigm for efficient and scalable ANN search. As the name suggests, HNSW searches a layered hierarchical graph to quickly identify neighborhoods of similar points to a given query vector. But is this hierarchy even necessary? A rigorous experimental analysis to answer this question would provide valuable insights into the nature of algorithm design for ANN search and motivate directions for future work in this increasingly crucial domain. To that end, we conduct an extensive benchmarking study covering more large-scale datasets than prior investigations of this question. We ultimately find that a flat navigable small world graph graph retains all of the benefits of HNSW on high-dimensional datasets, with latency and recall performance essentially *identical* to the original algorithm but with less memory overhead. Furthermore, we go a step further and study *why* the hierarchy of HNSW provides no benefit in high dimensions, hypothesizing that navigable small world graphs contain a well-connected, frequently traversed "highway" of hub nodes that maintain the same purported function as the hierarchical layers. We present compelling empirical evidence that the *Hub Highway Hypothesis* holds for real datasets and investigate the mechanisms by which the highway forms. The implications of this hypothesis may also provide future research directions in developing enhancements to graph-based ANN search.

## 1 INTRODUCTION

Near neighbor search is a fundamental problem in computational geometry that lies at the heart of countless practical applications. From industrial-scale recommendation [14] to retrieval-augmented generation [27] and even to computational biology [50], numerous data-intensive tasks utilize similarity search at some location in the stack. As a result, similarity indexes are very well-studied [1, 17, 23, 25, 32] with multiple large-scale benchmarks and leaderboards to compare techniques [3, 42].

Historically, the state-of-the-art for near neighbor search involved constructing sophisticated tree-based data structures, such as $kd$-trees [4] and cover trees [6], that guaranteed exact solutions while avoiding a brute-force examination of all points. However, the recent advent of large-scale neural representation learning, including large language models (LLMs), places a significant strain on these classical methods that were developed to target a much lower-dimensional search space. In response, the community has turned to approximate search methods. While alternative approximate indexing methods such as locality-sensitive hashing [18] and

product quantization [24], have garnered significant interest, graph-based approaches generally achieve the strongest performance on established ANN benchmarks [3, 42]. Introduced in 2016, the Hierarchical Navigable Small World (HNSW) algorithm [32], emerged as one of the first high-performance graph-based search indexes at scale and still enjoys immense popularity to this day with over 4300 Github stars[1].

As the name implies, a core feature of the HNSW index is its hierarchically layered graph akin to a skip list [38] where the search process iteratively traverses through graphs of increasing density before converging to a neighborhood of similar points in the final graph layer. By drawing intuition from skip lists, the HNSW authors argue that the initial coarse graph layers allow for efficiently identifying the neighborhood of similar points in the collection through fewer overall comparisons.

Despite the immense popularity of HNSW, however, the algorithm suffers from multiple scalability bottlenecks. For instance, the hierarchical layers introduce a hefty memory overhead. Moreover, Malkov and Yashunin [32] also note in Section 6 of their paper that the hierarchy can reduce the overall system throughput in distributed settings when compared to a flat NSW graph.

While the high cost of this hierarchy in HNSW is traditionally justified by the low latency of graph-based indexes, we ask the question of whether the hierarchy is even necessary in the first place. A growing body of results suggest that the hierarchy may be a vestigial artifact left behind by the lower-dimensional problems of the past. For example, Lin and Zhao [28] published a set of experiments suggesting that the hierarchical component of HNSW is only beneficial for low-dimensional inputs ($d < 32$). Other recent works have observed similar problems with the hierarchy. For example, Coleman et al. [8] present an ablation study showing the same behavior in the appendix of their paper. This behavior seems to be an increasingly well-known "folklore" of the similarity search community, but there is currently an absence of a thorough investigation and analysis of this phenomenon.

An exhaustive benchmark comparison on the utility of the HNSW hierarchy would go a long way towards improving our understanding of graph-based similarity search in high-dimensional space and inform future research directions in developing improvements. However, there are several key difficulties with executing such a benchmarking study. The first issue is that running a meaningful ablation study requires first engineering a similarity search implementation over a flat NSW graph that reaches performance parity with HNSW and other state of the art approaches [12, 23, 43]. This is an increasingly nontrivial task, as the community has invested significant performance engineering efforts into these codebases such that a non-expert implementation has little hope of being competitive.
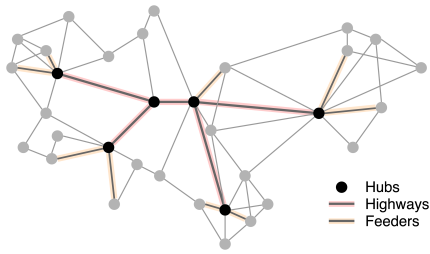
---

[1]https://github.com/nmslib/hnswlib

Blaise Munyampirwa, Vihan Lakshman, and Benjamin Coleman



**Figure 1: We hypothesize that in high dimensions, graph-based ANN indexes naturally form a "highway-feeder" structure, where a small subset of nodes and edges are easily reached, well-connected, and heavily traversed.**

The second challenge is designing benchmarking experiments in a manner that avoids the confounding effect of the performance of a particular software *implementation* when making conclusions about the efficacy of the *algorithm*. Finally, benchmarking near-neighbor search can be highly nuanced. For example, it is possible that hierarchical structures are strongly beneficial for isolated slices of the input. These may sharply reduce the $99^{\text{th}}$ percentile latency while not affecting the median.

Perhaps most importantly, we still have no satisfactory understanding of *why hierarchy does not help*. Hierarchical structures are a mainstay of algorithm design, where a common trick is to reduce an $O(n)$ search process to a sublinear one by traversing a (balanced) hierarchy [9, 16, 34, 38]. Arguably, it is counterintuitive for this idea to fail to hold in the context of high-dimensional similarity search – especially when we have strong positive results that hierarchy *helps* in low dimensions [6, 11, 28, 40]. Thus, an exhaustive benchmark and deeper analysis into the necessity of the hierarchy in HNSW would shed further light on the nature of algorithm design in high-dimensional spaces and thus may be of independent interest to the community as well.

## 1.1 Contributions

In this paper, we address the question of whether the hierarchical component of HNSW is truly necessary. Our central research question is, *"Can we achieve the same performance on large-scale benchmarks with simply a flat navigable small world graph?"* In summary, we make the following contributions:

**Benchmarking the hierarchy:** We rigorously benchmark HNSW to understand whether the hierarchy is necessary. To do so, we reproduce and extend the hierarchy ablations of previous studies, finding that, on high-dimensional vector datasets, it is indeed beneficial to remove the 'H' from HNSW.

**Why does hierarchy not help?** We hypothesize that the hierarchy benefits decrease in high-dimensions due to *hubness*. Hubness is a high-dimensional phenomenon that causes a skewed distribution in the near-neighbor lists of search queries [39]. We hypothesize that hubness leads to preferential attachment in the similarity search graph, inducing the formation of *easily-traversed highways* that connect disparate regions of the graph. This hypothesis, which we call the *Hub Highway Hypothesis*, explains why we

no longer need the hierarchy in high dimensions; we can simply traverse the intrinsic highway structure that naturally forms in high-dimensional spaces. To that end, we conduct a series of experiments to investigate the hub-highway hypothesis. Our results ultimately support this hypothesis, showing that hubness is responsible for driving the connectivity of similarity search graphs. This insight opens up exciting new research directions in graph construction, link pruning, and graph traversal.

Our specific contributions are as follows.

- We release an implementation for a flattened version of HNSW[2], called FlatNav[3], that reaches performance parity with the original version with considerable memory savings.
- We demonstrate that hierarchy does not improve performance in either the median or tail latency case by building HNSW and FlatNav indexes over 13 popular benchmark datasets ranging in size from 1 million to 100 million vectors.
- We conduct an analysis of hubness phenomena in high-dimensional metric spaces and the resulting HNSW graphs, finding strong empirical support for the hub-highway hypothesis.

**Practical implications:** Our benchmarks reveal that HNSW can be significantly optimized for modern high-dimensional embedding workloads. For instance, our implementation saves roughly 38% and 39% of peak memory consumption on two Big-ANN benchmark datasets compared to `hnswlib` (and sizable further headroom is likely). Our results confirm the folklore of the similarity search community, conclusively demonstrating that we can remove the hierarchy on high-dimensional inputs with impunity.

## 2 BACKGROUND: SIMILARITY SEARCH & HNSW

In the similarity search (or $k$-NNS) problem, we are interested in retrieving $k$ elements from a dataset $\mathcal{D} = \{\mathbf{x}_i, \ldots, \mathbf{x}_n\} \subset \mathbb{R}^d$ that minimize the distance to a given query $q \in \mathbb{R}^d$ (or, equivalently, maximize the vector similarity). More precisely, given a similarity function $\phi : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$, the nearest neighbor $\mathbf{x}^* \in \mathcal{X}$ of $q$ is defined as

$$\mathbf{x}^* := \underset{\mathbf{x}_i \in \mathcal{D}}{\arg\max}\, \phi(\mathbf{x}_i, q)$$

where $\phi$ is usually the $\ell_2$ or cosine similarity. With the enormity of modern data workloads and the underlying vector dimensionality, it becomes computationally infeasible to exhaustively search for the true top-$k$ neighbors for any query $q$. Thus, approximate search algorithms trade-off quality of the search for lower latency.

In the approximate nearest neighbor search (ANNS) regime, we evaluate the quality of the search procedure typically by the Recall@$k$ metric. More formally, suppose a given ANNS search algorithm outputs a subset $O \subseteq \mathcal{D}, |O| = k$, and let $G \subseteq \mathcal{D}$ be the true $k$ nearest neighbors of a query $q$. We define Recall@$k$ by $\frac{|O \cap G|}{k}$. ANNS algorithms seek to maximize this metric while retrieving results as quickly as possible.

---

[2]https://github.com/nmslib/hnswlib
[3]https://github.com/BlaiseMuhirwa/flatnav

## 2.1 HNSW Overview

With this formalization of the ANNS problem, we will now briefly review the key elements of the HNSW algorithm, which is the central focus of our benchmarking study. As we alluded to previously, HNSW builds off of prior work in *navigable small world graph* indexes introduced in [31]. Small world graphs are a well-studied phenomenon in both computing and the social sciences and are primarily defined by the fact that the average length of a shortest path between two vertices is small (typically scaling logarithmically with the number of nodes in the network) [26, 46, 48]. Small world graphs are also often characterized by the presence of well-connected *hub nodes* which we discuss further in the next section.

---

**Algorithm 1** HNSW Construction

1: **Input:** Set of data points $D$, max layer $L_{max}$, max connections per layer $M$, layer insertion probability $m_l$, size of dynamic candidate list $efc$
2: **Output:** HNSW graph with hierarchical layers
3: **procedure** CONSTRUCT($D, L_{max}, M, m_l, efc$)
4:     Initialize empty hierarchical graph $G$
5:     Initialize entry point $ep \leftarrow$ None
6:     **for** each $p \in D$ **do**
7:         $L_p \leftarrow$ GeometricDistribution($m_l$)
8:         **if** $ep =$ None **then**
9:             Set $p$ as entry point $ep$
10:            Insert $p$ into all levels $\leq L_p$
11:         **end if**
12:         **for** $l = L_{max}$ to $L_p$ **do**
13:             $ep \leftarrow$ SearchLayer($G, l, p, ep, efc$)    ▷ Algorithm 2
14:         **end for**
15:         **for** $l = 0$ to $L_p$ **do**
16:             $N \leftarrow$ SelectNeighbors($p, G, l, M$)
17:             Add edges from $q$ to each neighbor $n \in N$ at layer $l$
18:             **if** $n \in N$ has $< M$ edges **then**
19:                 Add back-connections to $q$ to node $n$.
20:             **else**
21:                 Run SelectNeighbors on $\{q,$ edges of $n\}$.
22:             **end if**
23:         **end for**
24:     **end for**
25:     **if** $L_p > L_{ep}$ **then**
26:         $ep \leftarrow p$
27:     **end if**
28: **end procedure**
29:
30: **function** SELECTNEIGHBORS($p, G, l, M$)
31:     Compute distances from $p$ to all nodes in $G[l]$
32:     Return $M$ nodes based on selection heuristic in [2]
33: **end function**

---

## 2.2 HNSW Search Algorithm

While small world graphs are, by construction, suited for efficient greedy graph traversal, the HNSW authors argue that the polylogarithmic scaling of the search process is still too inefficient for

---

**Algorithm 2** HNSW Query

1: **Input:** Graph $G$, layer $l$, query $q$, starting point $p$, number of nearest neighbors to return $efs$
2: **procedure** SEARCHLAYER($G, l, q, p, efs$)
3:     Candidate queue $C = p$, currently top results queue $T = p$, visited list $V = p$
4:     **while** $C$ is not empty **do**
5:         $c \leftarrow$ nearest element from $C$ to $q$
6:         $f \leftarrow$ furthest element from $T$ to $q$
7:         **if** dist(c, q) > dist(f, q) **then return** T
8:         **end if**
9:         **for** $e \in$ neighbourhood($c$) at layer $l$ **do**
10:             **if** $e \in V$ **then**
11:                 continue
12:             **end if**
13:             $V.add(e)$
14:             **if** $dist(e, q) \leq dist(f, q)$ or $|T| \leq efs$ **then**
15:                 $C.add(e)$
16:                 $T.add(e)$
17:             **end if**
18:             **if** $|T| \geq efs$ **then**
19:                 Remove furthest point to q from T
20:             **end if**
21:             $f \leftarrow$ furthest element from $T$ to $q$
22:         **end for**
        **return** T
23:     **end while**
24: **end procedure**

---

the demands of near neighbor search on large datasets. This claim motivates the design of HNSW where the hierarchy allows for computing a fixed number of distances in each graph layer independent of the network size.

Specifically, the HNSW index is constructed in an iterative fashion. For a newly inserted element $x$, the algorithm will randomly select a maximum layer $l$ and then insert the new point into every layer up to $l$. This randomized process is executed with an exponentially decaying probability distribution such that, in expectation, each subsequent layer has exponentially more nodes than its predecessor. Within a layer, HNSW greedily adds edges between $x$ and its $M$ closest neighbors (where $M$ is a hyperparameter) where the neighbors consist of previously inserted points. This process then repeats in the subsequent layer below using the closest neighbors found in the prior graph as entry points. Through this process, the top layer of the hierarchy will be the coarsest directed graph, consisting of the fewest nodes and edges, and the bottom layer will be the densest and contain all of the nodes, each with connections to (up to) $M$ neighbors. As an additional, and important, optimization, HNSW also implements the pruning heuristic of [2] that will prune an edge from $u$ to $v$ if there exists another edge from $u$ to a neighbor $w$ of $v$ such that the distance from $u$ to $w$ is less than that of $u$ to $v$.

The search procedure of HNSW, described in Algorithm 2 also executes iteratively where the algorithm maintains a list of candidate points at each layer of the hierarchical graph before returning
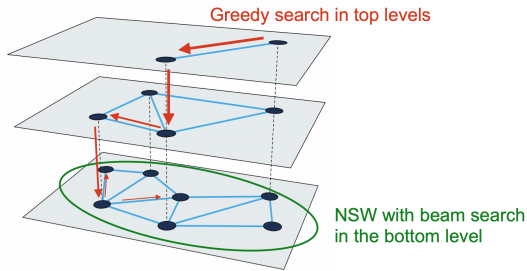
**Figure 2: Visualization of the HNSW search index. We benchmark the full hierarchical graph search procedure against simply searching in the bottom base layer.**

the final list $k$ nearest neighbor candidates after traversing the base graph layer.

## 3 RELATED WORK

### 3.1 Near-Neighbor Benchmarks

Due to the importance of the $k$-NNS problem, there have been several large-scale near-neighbor benchmarks in recent years. The original ANN Benchmarks [3] made significant impact because it provided the first standard evaluation of near-neigh3bor search algorithms. Since its inception, the benchmark has grown in scope from a handful of algorithms to include over 30 methods on 9 datasets. However, most of the datasets included in ANN Benchmarks are relatively small by modern standards at around one million points. To reflect the growing need for high-dimensional, large-scale embedding search, NeurIPS hosted a competition track known as Big ANN Benchmarks [42], which evaluates algorithms on five billion-scale datasets. Graph-based indexes easily and consistently place among the top three algorithms in each benchmark. However, these benchmarks are not sufficient for our purposes because they focus exclusively on overall throughput metrics and do not report tail latency numbers such as the 99[th] percentile, which is important in our case for teasing apart the impact of hierarchy in HNSW.

**Hierarchy Studies:** Hierarchical structures are considered a fundamental component of graph-based near-neighbor algorithms. As a result, numerous high-performing algorithms use some form of hierarchy including HNSW [32], ONNG [21] PANNG [20], and HC-NNG [35]. However, this practice has recently come under scrutiny. Dobson et al. [10] show that the hierarchy may be unnecessary for some workloads, based on the observation that HNSW underperforms both HCNNG (which uses a shallower hierarchy) and DiskANN [23] (which does not use a hierarchy). Lin and Zhao [28] present results on synthetic data showing that the hierarchy is only beneficial for low-dimensional data ($d < 32$), but analyze only a limited number of small real-world datasets and – most importantly – do not explain why the hierarchy fails to provide value. In this work, we aim to address these gaps in the literature through more comprehensive benchmarks coupled with an analysis of why hierarchical structures in high-dimensional space may not add value.

In the following section, we aim to reproduce the results of both Malkov and Yashunin [32] and Lin and Zhao [28] to confirm that we can independently replicate their findings with our own software implementation before proceeding to new experiments on larger-scale benchmark datasets.

### 3.2 Hubness in High Dimensional Spaces

Astute readers might observe that the HNSW graph construction algorithm described above does not explicitly enforce the small world property and instead adds edges between nodes based on their proximity in the metric space. The connection between proximity and the small world property arises due to *hubness*.

Specifically, hubness is a property of high-dimensional metric spaces where a small subset of points (the "hubs") occur a disproportionate number of times in the near-neighbor lists of other points in the dataset [39]. In other words, a small fraction of nodes are highly connected to other points.

A related, but slightly different, notion of hubness that has been at the center of theoretical analysis is concentration of distances in high-dimensional spaces. It has been shown that, in expectation, the $\ell_2$ distance between independent and identically distributed (i.i.d) vectors grows with $\sqrt{d}$ while the variance tends to a constant as $d$ approaches infinity [45]. As a result, the $\ell_2$ distance loses its discriminative power as $d$ increases. In fact, this concentration phenomenon is not only constrained to $\ell_2$ but also applies to other $\ell_p$ and even fractional norms [15].

Due to undesirable consequences of the hubness phenomenon, such as poor clustering quality, a large body of work has focused on hubness reduction strategies. For instance, [49] introduced local scaling which scales distances $d(\mathbf{x}, \mathbf{y})$ by accounting for local neighborhood information. Interestingly, our work stands in contrast to this literature on hubness reduction by presenting a case study where hubs provide tangible value in an algorithmic setting, namely in accelerating greedy traversal in near neighbor proximity graphs. This result may be of independent interest to machine learning and algorithms researchers as well.

## 4 REPRODUCTION OF PRIOR STUDIES

In this section, we present a replicability study using four `flatnav` NSW implementation. In particular, we revisit the experimental design of two prior works in the literature: the original 2016 HNSW paper of Malkov and Yashunin [32] and a subsequent 2019 paper from Lin and Zhao [28] that found limitations with the hierarchical component of HNSW. As we discussed in the previous section, these prior works possess limitations in experimental design, scope of benchmarking datasets, and a lack of analysis into understanding the results, which motivates our work in this paper. Nevertheless, we use these prior studies as a starting point to see if we can independently replicate these results via our own FlatNav implementation. Such a reproduction would both further validate the soundness of these previous experiments over the test of time as well as provide confirmation of the correctness of FlatNav before we proceed to new, larger-scale benchmarks.

Following the same setups as Malkov and Yashunin [32] and Lin and Zhao [28] we generate a series of random vector datasets of varying dimensionality where each vector component is sampled

uniformly at random from the range $[0, 1)$. In particular, we consider dimensionalities of $d = 4, 8, 16$ and $32$. As in [28], we set the number of near neighbors to retrieve to $k = 1$ (departing from the default of $k = 100$ we use elsewhere in this paper). We also tried including the `sw-graph` NSW baseline [7] that Malkov and Yashunin [32] used in their evaluation to benchmark against HNSW, but we were not able to run this older library successfully. However, we were able to replicate these prior findings using our own `flatnav` implementation which is conceptually identical to `sw-graph` but with more software optimizations to achieve engineering parity with `hnswlib` (detailed in Section 5.1).
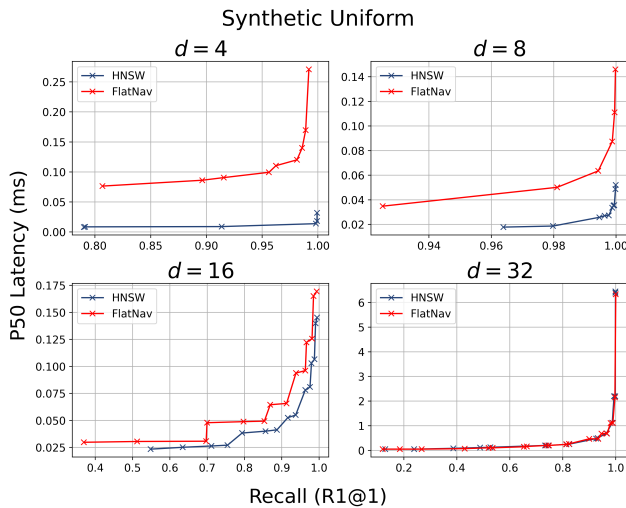


**Figure 3: Median Latency vs. Recall of HNSW and FlatNav across dimensions** $d = 4, 8, 16, 32$**. We observe that the hierarchical structure accelerates search only when** $d < 32$**, matching the findings of Lin and Zhao [28]. Our results demonstrating a significant advantage with HNSW on synthetic datasets with dimensionality** $d = 4$ **and** $d = 8$ **also match the findings of the original HNSW paper Malkov and Yashunin [32]**

.

As shown in Figure 3, we successfully replicated the experiments benchmarking HNSW versus a flat NSW graph from two prior research papers. Notably, both of these previous works primarily experiment with randomly generated vector data with very low dimensionality by the standards of modern machine learning. Coupled with our findings in the next section where we find no discernible difference between HSNW and a flat graph index on high-dimensional datasets, our results suggest a simple decision criterion for selecting a search index: **For dimensionality** $d < 32$**, HNSW and the hierarchy provide a speedup. Otherwise, the simplicity and memory savings of a flat NSW index provide more benefit.**

We also had the opportunity to discuss our findings with the lead author of [32] who confirmed that the hierarchy provides a robust speedup on these low-dimensional datasets but noted the performance on higher dimensional vectors remained less clear, which further motivated us to take up the benchmarking study in the next section.

# 5 FLATNAV BENCHMARKING EXPERIMENTS

In this section, we report the results of our benchmarking study comparing the performance of flat HNSW search to hierarchical search on a suite of standard high-dimensional benchmark datasets drawn from real machine learning models.

A challenge in designing such benchmarking experiments is separating the confounding impact of performance engineering in the implementation of an approach from any purely algorithmic advantages of a method. We argue that this is a pitfall that has clouded several past benchmarking studies in ANN search. To circumvent this challenge, we fix the implementation in our experimental design such that the *same code* is used in constructing the base navigable small world graph. In particular, we use the popular `hnswlib` library as our baseline HNSW implementation. To benchmark the flat NSW index performance, we extract the bottom graph layer from `hnswlib` and reimplement the HNSW's greedy search heuristic over this flat graph. To the best of our knowledge, we provide the first benchmarking study on the utility of HNSW's hierarchy that makes an effort to achieve parity in the implementations of the two techniques to center the focus on the central algorithmic question.

## 5.1 Software Performance Optimizations

One major challenge in prior near neighbor benchmarking studies is the confounding effect of software performance optimizations when making conclusions about algorithmic efficacy. To address this issue, our implementation of a flat similarity search graph in `flatnav` closely follows the performance optimizations native to `hnswlib`. We use SIMD vector instructions to efficiently compute the distance between any two vectors, which is the key operation underpinning similarity search. Depending on the hardware type, we provide distance computations using SSE, AVX and AVX-512 intrinsics for `float`, `uint8_t` and `int8_t` data types. The choice of which intrinsics family to use is a compile-time choice, which is determined based on the dataset dimensionality and hardware support, similar to `hnswlib`. For all benchmark experiments, we use the `float` data type.

Furthermore, one common problem in constructing and querying graph-based indexes is the random data access pattern which can hamper latency and throughput. To improve search throughput, `hnswlib` implements prefetching using the `_mm_prefetch` intrinsic which reduces memory access latency by prefetching data into the L1 cache. Since caching can significantly affect latency metrics, we control for this performance optimization by leveraging the same intrinsic in the `flatnav` implementation so that our benchmarks strictly adhere to algorithmic differences.

## 5.2 Datasets and Compute

We utilize the benchmark datasets released through the popular leaderboards ANN Benchmarks [3] and Big ANN Benchmarks [42]. The specific datasets and their associated statistics are presented in Table 1. For the Big ANN Benchmark datasets, we consider both the 10M and 100M collection of vectors for which the ground truth near neighbors have previously been computed and released. We did not experiment with the largest Big ANN datasets with 1 billion vectors since constructing HNSW indexes at this scale requires over

1.5TB of RAM, which exceeded our compute resources. Below, we include our benchmarking results for the four 100M-scale datasets available through Big ANN Benchmarks. We see that our flat HNSW implementation achieves performance parity with the hierarchical HNSW implementation.

| Dataset | Dimensionality | # Points | # Queries |
|---------|----------------|----------|-----------|
| BigANN[†] | 128 | 100M | 10K |
| Microsoft SpaceV[†] | 100 | 100M | 29.3K |
| Yandex DEEP[†] | 96 | 100M | 10K |
| Yandex Text-to-Image[†] | 200 | 100M | 100K |
| GloVe | {25, 50, 100, 200} | 1.2M | 10K |
| NYTimes | 256 | 290K | 10K |
| GIST | 960 | 1M | 1K |
| SIFT | 128 | 1M | 10K |
| MNIST | 784 | 60K | 10K |
| DEEP1B | 96 | 10M | 10K |

**Table 1: Dataset Statistics. The datasets marked by † are from the BigANN benchmarks [42]. The remaining are taken from ANN Benchmarks [3].**

For our benchmarks on datasets consisting of fewer than 100M vectors in the collection, we use an AWS c6i.8xlarge instance with an Intel Ice Lake processor and 64GB of RAM. We selected this particular public cloud instance to facilitate accessible reproducibility of our experiments. For the 100M-sized large-scale experiments, we use a cloud server equipped with an AMD EPYC 9J14 96-Core Processor and 1 TB of RAM.

## 5.3 Latency Results

*5.3.1 BigANN Benchmarks [42].* In Figures 4 and 5, we compare latency metrics for HNSW and FlatNav at the 50th and 99th percentile for the four 100M datasets from BigANN benchmarks listed in Table 1. All of our results support the conclusion that `flatnav` achieves nearly identical performance to `hnswlib`.

From the results in Figures 4 and 5, we observe that there is no consistent and discernable gap between FlatNav and HNSW in both the median and tail latency cases. These results suggest that the hierarchical structure of HNSW provides no tangible benefit on practical high-dimensional embedding datasets.

*5.3.2 ANN Benchmarks [3].* Moreover, we repeat the same experimental setup comparing HNSW and FlatNav on the ANN Benchmark datasets listed in Table 1. In Figures 6 and 7, we report the p50 and p99 latency of all of the non-GloVe ANN Benchmarks. Although these datasets are smaller in scale than the BigANN Benchmarks, we still see no discernible difference in latency between HNSW and FlatNav which supports our hypothesis that the vector dimensionality and not the size of the collection is the main driver of eliminating the need for hierarchical search in small world graphs. We see further evidence of this idea in Table 8, where we plot the median latency versus recall on the four GloVe benchmark datasets ranging in dimensionality from 25 to 200. Ultimately, these results, coupled with our reproducibility studies in the previous section, provide compelling evidence that for high-dimensional datasets, which are the standard in modern similarity search workloads, there is no apparent performance benefit to the hierarchical layers of HNSW.
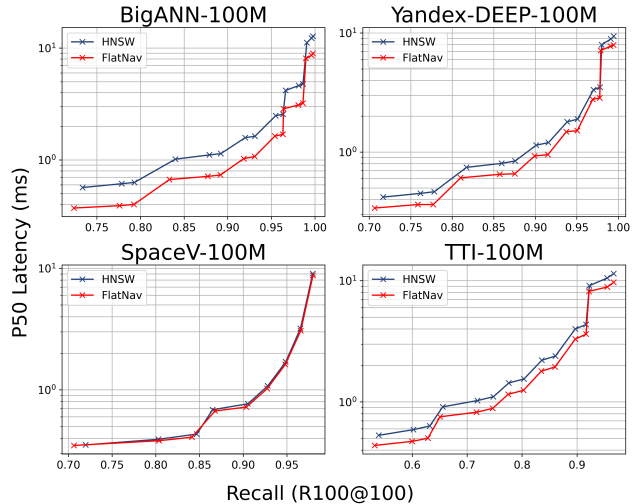


**Figure 4: p50 Latency vs. Recall for HNSW and FlatNav over four 100M-scale benchmark datasets with dimensionality between 96 and 200. We observe that the flat NSW index achieves essentially identical performance to the hierarchical HNSW index.**
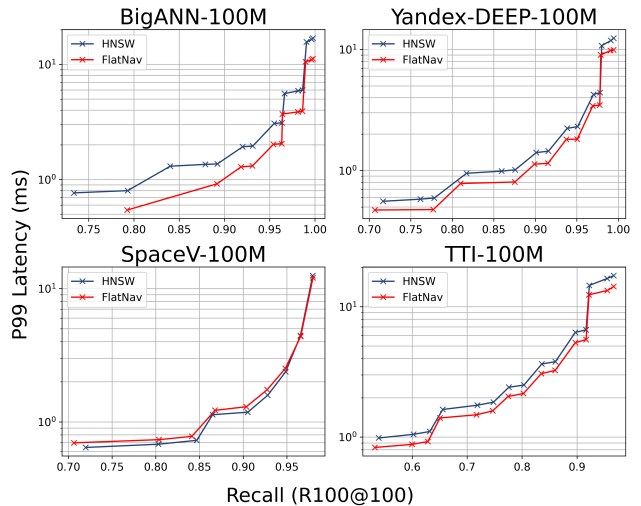


**Figure 5: p99 Latency vs. Recall for HNSW and FlatNav over four 100M-scale benchmark datasets with dimensionality between 96 and 200. As with the median case, we observe nearly identical performance between the two indexes.**

## 5.4 Memory

In this section we measure the memory savings from removing the hierarchy by running a memory profiler during index construction. In terms of memory allocation, similarly to `flatnav`, `hnswlib` allocates static memory during index construction comprising base layer node allocation, a visited node list and a list of mutexes in the
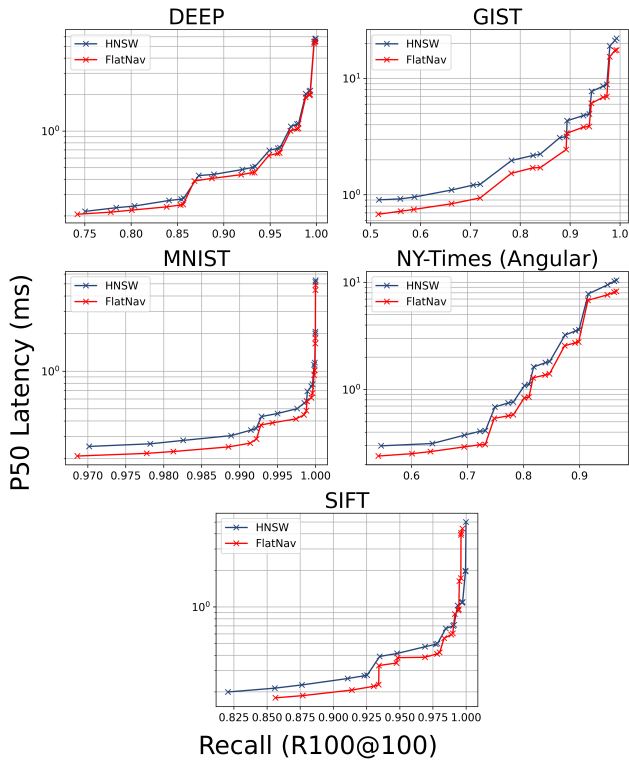
Figure 6: p50 Latency vs. Recall for HNSW and FlatNav over five ANN Benchmark datasets ranging in dimensionality from 96 to 960. The flat NSW index again achieves practically indistinguishable performance from the hierarchical graph index.

| Dataset | # Data | hnswlib Memory | flatnav Memory |
|---------|--------|----------------|----------------|
| BigANN | 100M | 183 | 113 |
| Microsoft SpaceV | 100M | 104 | 85.5 |
| Yandex DEEP | 100M | 100 | 60.7 |

Table 2: Peak Index Construction Memory in GBs. We observe that flatnav requires considerably less memory during construction compared to hnswlib.

multi-threaded setting. Additionally, it also incurs memory cost attributable to the hierarchy, particularly maintaining the dynamically allocated links between nodes at each layer.

We benchmarked both libraries against a subset of the BigANN benchmarks. Table 2 shows the peak memory allocated by the two implementations during index construction for the BigANN, Microsoft SpaceV and Yandex DEEP benchmarks. Since multithreading has a runtime overhead, we fix the number of cores to 32 in each one of the stated benchmark. For BigANN, we observe a 38% reduction in peak memory, a 39% reduction for Yandex DEEP, and an 18% reduction for the Microsoft SpaceV benchmark. This shows that we are able to save significant memory by removing the hierarchy, and it is likely that we can optimize flatnav implementation to save memory further.
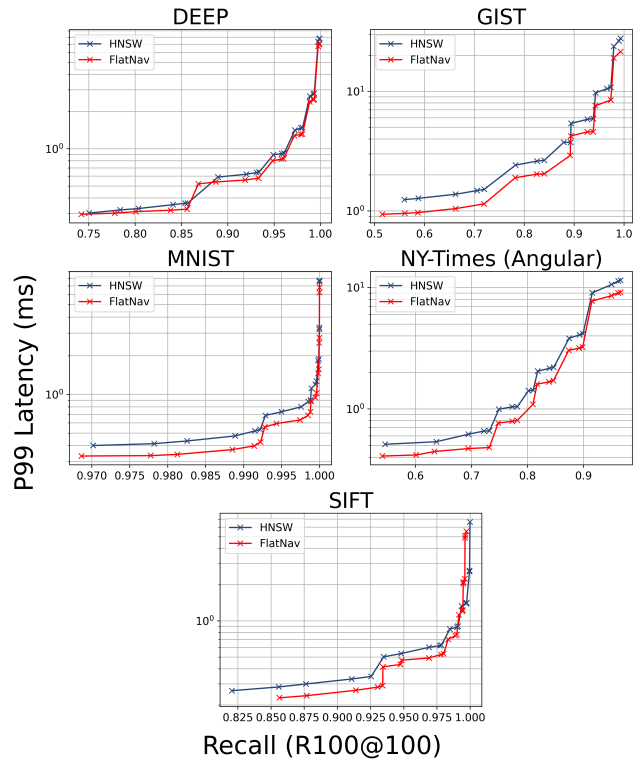
Figure 7: p99 Latency vs. Recall for HNSW and FlatNav over five ANN Benchmark datasets ranging in dimensionality from 96 to 960.
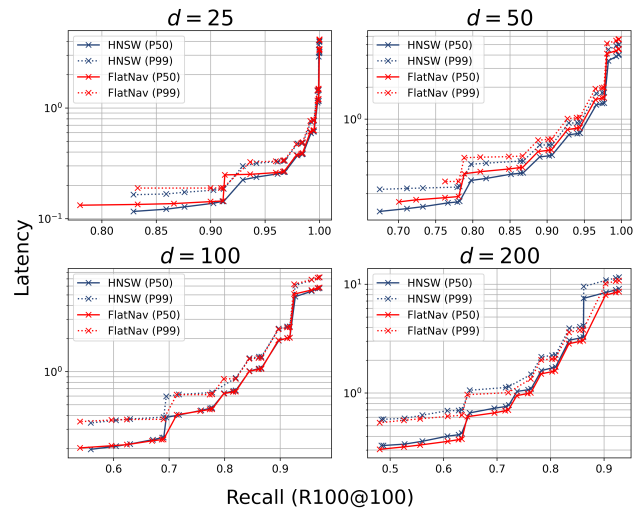


Figure 8: p50 and p99 Latency vs. Recall for HNSW and FlatNav over the GloVe ANN Benchmark datasets with respective dimensionalities 25, 50, 100, and 200. We observe that HNSW has a minor advantage in low dimensions and that FlatNav has a minor advantage in high dimensions, but also that there is no clear difference in performance.

One caveat of these reported memory savings is that we are comparing different two software implementations in `hnswlib` and `flatnav`. Since `hnswlib` is a mature library widely used by practitioners as well as researchers, it supports more features than `flatnav` and thus must maintain additional complexity whereas our implementation, while performant, is more of a research prototype. Therefore, differences in code may account for a significant part of the peak memory usage differences. Nevertheless, we believe our findings are relevant and still noteworthy given that `hnswlib` is so widely adopted. By demonstrating that we can considerably reduce the memory overhead of `hnswlib` without sacrificing performance, we hope to bring the community's attention to the opportunities for further optimization in this direction.

## 6 THE HUB HIGHWAY HYPOTHESIS

In the previous section, we established through a series of careful benchmarking experiments that there appears to be no discernible difference in performance between HNSW and its flat counterpart with no hierarchical graph structure. We now turn our attention to studying why the hierarchy appears to provide no benefit in the search process.

In our experiments, we observed that a small fraction of nodes appear in the set of near neighbors for a disproportionate number of other vectors. We thus conjecture that the hub structure prevalent in high-dimensional data performs the same functional role as the hierarchy.

HYPOTHESIS (HUB HIGHWAYS). *In high-dimensional metric spaces, $k$-NN proximity graphs form a highway routing structure where a small subset of nodes are well-connected and heavily traversed, particularly in the early stages of graph search.*

We remark that the existence of hub nodes in high-dimensional space is not a new observation [39]. The novelty of our hypothesis lies in connecting the idea of hubness to the notion of accelerating near neighbor search in ANN proximity graphs. In particular, we conjecture that near neighbor queries over proximity graphs in high dimensions often spend the majority of their time visiting hub nodes early on in the search process before converging to a local neighborhood of near neighbors. This procedure succeeds because hub nodes are very well connected to other parts of the graph and thereby efficiently route queries to the appropriate neighborhood in much the same manner that the layered hierarchy purports to do.

In many applications, such as clustering, hubness is considered an undesirable property that harms algorithmic performance, motivating the need for hubness reduction techniques [49]. On the contrary, our *Hub Highway Hypothesis* presents a case where hubness offers a distinct advantage in a practical algorithmic setting.

In the remainder of this section, we present an experimental design and a series of results that provide empirical evidence in the affirmative for the existence of such a highway routing mechanism amongst hubs in navigable small world graphs.

### 6.1 Intuition:

The formation of a hub highway can be understood through concentration of measure phenomena.

**Concentration of distances:** Under many data-generating distributions, the expected distances between randomly-drawn points converge to a constant value as dimension increases [5]. This suggests that in sufficiently high-dimensional spaces, we do not need to explicitly encourage the formation of "long-range" connections. Since all distances are approximately equal, all connections have similar length and a few edge traversals should span the graph of a high-dimensional dataset.

**Concentration of measure:** A well-known fact about high-dimensional geometry is that volume concentrates at the extrema of a shape (e.g., the surface of a ball or the faces of a hypercube). These portions of the distribution contain a disproportionate number of points and frequently contain the hubs of the metric space. For example, Low et. al. demonstrate that the the hubs of a uniform hypercube distribution are located at the corners [29], which are exactly the points we would wish to connect via long-range connections in a hierarchical graph.

### 6.2 Methodology

**Argument sketch:** We will demonstrate the *Hub Highway Hypothesis* by providing empirical evidence for the following claims.

(1) Some nodes are visited by queries much more frequently than others. The relative popularity of these *hub nodes* is explained by the hubness phenomena that arises in high-dimensional datasets.

(2) The hub nodes tend to be connected to each other, forming a well-connected subgraph of hubs.

(3) Queries visit many hub nodes early in the search process, before visiting less well-traversed neighborhoods.

**Empirical measures of hubness:** To support the first part of our argument, we will need a formal definition and characterization of hubness.

Following [39], let $\mathbf{x}, \mathbf{x}_1, \ldots, \mathbf{x}_n$ be a collection of vectors drawn from the same probability distribution with support $\mathcal{S} \subseteq \mathbb{R}^d$, and let $\phi : \mathcal{S} \times \mathcal{S} \to \mathbb{R}$ be a distance function.

Furthermore, for $1 \leq i, k \leq n$, let $p_{i,k}$ be defined by

$$p_{i,k}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \text{ is among the } k\text{-NN set of } \mathbf{x}_i \text{ under } \phi \\ 0 & \text{otherwise} \end{cases}$$

Now let $N(\mathbf{x})$ be the random variable defined by

$$N(\mathbf{x}) := \sum_{i=1}^{n} p_{i,k}(\mathbf{x})$$

which represents the number of vectors that have $\mathbf{x}$ included in their $k$-nearest neighbors.

For any dataset $\mathcal{D}$, we can compute $N(\mathbf{x}), \mathbf{x} \in \mathcal{D}$, which yields a discrete distribution $N_k$. We are interested in the skewness of this distribution, given by

$$S_{N_k} = \frac{\mathbb{E}\left[\left(N_k - \mu_{N_k}\right)^3\right]}{\sigma_{N_k}^3}$$

This measure characterizes the asymmetry of the $k$-occurrence distribution $N_k$, and it is the metric most often used to estimate the presence of hubs. Thus, in our experiments we use $S_{N_k}$ to quantify the presence of hubness in a particular dataset. The more skewed

the distribution of $N_k$, the more dataset $\mathcal{D}$ tends to have clusters of vectors (hubs) that keep regularly occurring in the $k$-NN sets of other vectors.

**Table 3: Datasets**

| Dataset | Dimensionality | # Data | # Queries |
|---|---|---|---|
| GIST | 960 | 1M | 1k |
| GloVe | 100 | 1.2M | 10k |
| NYTimes | 256 | 290K | 10k |
| Yandex-DEEP | 96 | 10M | 10k |
| Microsoft-SpaceV | 100 | 10M | 29.3k |
| IID Normal | {16, 32, 64, 128, 256, 1024, 1536} | 1M | 10k |
| IID Normal | {16, 32, 64, 128, 256, 1024, 1536} | 1M | 10k |

**Synthetic and ANN Benchmark Datasets**: To illustrate evidence for the Hub Highway Hypothesis, we experiment with both real and synthetic datasets as shown in Table 3. In addition to a subset of ANN Benchmark datasets, we generate synthetic datasets by drawing vectors from the standard normal distribution. The choice of the normal distribution is important because, as [29] observe, data sampled from the normal distribution exhibit a stronger hubness phenomenon due to strong density gradient than other distributions, such as sampling uniformly from a $d$-dimensional hypercube. Furthermore, since the hubness phenomenon is accentuated in high-dimensional spaces, we generate random data with increasing dimensionality which allows us to study the effect of these intrinsic highway structures that form as $d$ increases.

In the following three experiments we show empirical evidence supporting the Hub Highway Hypothesis. We start by examining the distribution of the number of times each node in the similarity search index is visited during search given a fixed number of queries, which we denote by $P_m(\mathbf{x}_i)$ to indicate that node $\mathbf{x}_i$ is visited exactly $m$ times during a fixed number of $k$-NN queries for each one of the datasets in Table 3. We show that this distribution is skewed to the right, thus confirming that certain nodes (*highway nodes*) are visited a disproportionate number of times. Next, using this distribution, we selectively choose the most visited nodes in the similarity search graph to be the hub node clusters and show evidence that these nodes are more connected to each other than random nodes using hypothesis tests. In the last experiment, we show that not only are these highway nodes more connected, they also allow for faster graph traversals for queries, hence supporting the claim that we no longer need hierarchy in high dimensional vector search where fast query times evolve as a result of the presence of the hub-highways.

## 6.3 Skewness of the Node Access Distribution

As described above, for this study we consider the discrete distribution $P_m(\mathbf{x}_i)$, $\mathbf{x}_i \in \mathcal{D}$, where $\mathcal{D}$ is one of the datasets in Table 3. For all these experiments we construct a HNSW index with parameters listed in Table 4 and import the base layer of the graph into FlatNav to construct a flat NSW index.

Figure 9 shows the log-normalized node access count distribution for different datasets. We observe that as the dimension $d$ increases, this distribution becomes more skewed to the right for $\ell_2$ distance-based datasets, indicating that a subset of nodes are visited much more often than the rest of the nodes in the graph. The cosine

**Table 4: Similarity search index parameters**

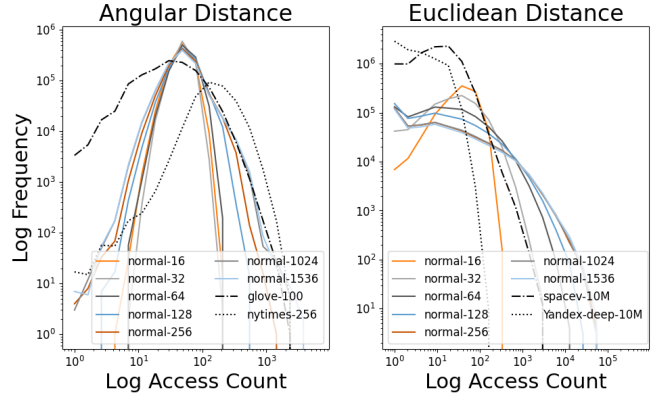| $m$ | $ef$-construction | $ef$-search | $k$ |
|---|---|---|---|
| 32 | 100 | 200 | 100 |



**Figure 9: Log-normalized node access count distribution** $P_m(\mathbf{x}_i)$ **for datasets using angular distance (left) and $\ell_2$ distance (right). Note that the x-axis is also log-normalized, so all histograms are skewed to the left – but the $\ell_2$ datasets have a much greater skew.**

distance, on the other hand, is known to have anti-hub properties that prevent such a dramatic skew even for $d \in \{1024, 1536\}$. This is consistent with the results observed by [39]. The increased skewness of the $P_m(\mathbf{x}_i)$ distribution as $d$ increases suggests that highway nodes occur as a result of the intrinsic dimensionality of the datasets.

**Confounding mechanisms:** While Figure 9 confirms the existence of hub nodes in the dataset, it does not distinguish between hubs that arise from the properties of the underlying metric space and hubs that form through some other mechanism. It is possible that *preferential attachment* explains the formation of hubs, since NSW graphs are built incrementally by sequentially adding points to an existing graph. Nodes that are added early in graph construction may become hubs by accumulating a greater-than-average share of inbound graph links, rather than by being popular neighbors in the metric space.

To investigate the effects of preferential attachment, we computed the variance ($R^2$ of a linear model) in the empirical node access distribution explained by the insertion ordering(Table 5). We log-transformed both the node access count and the insertion order before running the linear model and confirmed that the residuals are approximately normal by visually examining the QQ plots ($p < 10^{-6}$ for all models).

We observe a modest effect from the node insertion order in our synthetic data. This is particularly true for the angular datasets, which we believe to be due to the weaker hubness phenomena produced by the angular distance metric. Preferential attachment may account for a relatively greater share of the node access distribution when metric hubs are not present to heavily skew the distribution.

Ideally, we would repeat this analysis using the $K$-occurrence distribution, to show that the hubness of the metric space is more strongly predictive of the node acccess count than the insertion order. Unfortunately, it is not feasible to compute the $K$-occurrence distribution due to the $O(n^2)$ brute-force computation cost. However, we believe that the results in Table 5 still support the idea that the dimensionality of the metric space strongly contributes to the formation of hubs in the *Hub-Highway Hypothesis*, especially when combined with the evidence in Figure 9.

**Table 5: Correlation analysis for node insertion order and node access count, to determine effects of preferential attachment.**

| Dataset | Dimension | Explained Variance (%) |
|---|---|---|
| IID Normal (Angular) | 16 | 3.0% |
| IID Normal (Angular) | 32 | 8.7% |
| IID Normal (Angular) | 64 | 16.6% |
| IID Normal (Angular) | 128 | 23.3% |
| IID Normal (Angular) | 256 | 24.6% |
| IID Normal (Angular) | 1024 | 24.1% |
| IID Normal (Angular) | 1536 | 23.9% |
| IID Normal (L2) | 16 | < 0.1% |
| IID Normal (L2) | 32 | 3.1% |
| IID Normal (L2) | 64 | 7.6% |
| IID Normal (L2) | 128 | 8.7% |
| IID Normal (L2) | 256 | 7.1% |
| IID Normal (L2) | 1024 | 7.1% |
| IID Normal (L2) | 1536 | 6.6% |
| GloVe (Angular) | 100 | 0.2% |
| NYTimes (Angular) | 256 | 0.3% |
| GIST (L2) | 960 | < 0.1% |
| Yandex-DEEP (L2) | 96 | 0.3% |
| Microsoft-SpaceV (L2) | 100 | 0.5% |

## 6.4 Subgraph Connectivity of the Hub-Highway Nodes

Having demonstrated the occurrence of highway nodes, in the next set of experiments we present empirical evidence confirming that these nodes exhibit strong connectivity in the graph. One reasonable approach could be to identify hub nodes by thresholding the node access distribution (e.g., to the top 5% by access frequency) and determining whether these nodes form a connected component. However, this is strongly influenced by the threshold probability, and in any case the hubs need not form a single subgraph to support the *Hub Highway Hypothesis*.

Instead, we rely on statistical hypothesis testing to determine whether hub nodes are likely to connect to other hub nodes. We begin by explaining our procedure to identify hubs. Let $x_i, 1 \leq i \leq n$, be the vectors in a similarity search index for any dataset $\mathcal{D}$ under consideration.

- Using the empirical node access distribution $P_m(x_i)$, we identify hub nodes as those that fall into either the 95[th] or 99[th] percentile of node access counts. We assign a binary label to each node to indicate whether it is a hub. Let $h$ :

$\mathcal{D} \rightarrow \{0, 1\}$ be this assignment function with $h(x_i) = 1$ for nodes identified as hubs.
- We wish to estimate the likelihood that a randomly-chosen out-neighbor of a hub node is, itself, a hub. To do so, we examine the 1-hop out-degree expansion of each hub and count the number of adjacent hub nodes. By repeating this process for each hub node, we obtain a discrete distribution for the number of hubs to which each hub node is connected.
- Similarly, we select a set of random non-hub nodes from $V := \mathcal{D} \setminus \bigcup_{i=1}^{n} h(x_i) = 1$. For each node $x_i \in V$, we compute the same quantity to find the number of hubs with which non-hubs are connected, allowing us to construct the equivalent distribution for non-hubs.

Our hypothesis test investigates whether these two distributions differ. The null hypothesis is that hubs and non-hubs are equally likely to attach to hubs, and the alternative hypothesis is that hubs display different connectivity behavior than the rest of the graph. A finding that hubs are statistically more likely to attach to other hubs will support the *Hub Highway Hypothesis*.

We measure statistical significance using the standard two-sample parametric $t$-test as well as the non-parametric Mann-Whitney U-test [33], which is particularly suitable for ANN datasets because it does not require the normality assumption. Note that we consider dataset sizes larger than those often encountered in hypothesis testing ($n > 10^3$). A well-known consequence of the Central Limit Theorem (CLT) [13] is that the U-statistic used in non-parametric tests such as the Mann-Whitney test will converge to a normal distribution as the sample size increases. This, in turn, leads to $p$-values that approach 0 whenever there is a true difference between any two populations under consideration. For this reason, we also consider the effect size in our test and constrain the sample size to be 1000 for both populations.

Table 6 shows the results when we use the top-5% threshold to select hub nodes. Assuming the regular 0.05 threshold for rejecting the null hypothesis, note that in all by five cases we reject the null hypothesis under the Mann-Whitney $U$-test (and all by six cases using the two-sample $t$-test). Unsurprisingly, the effect size is much higher in the synthetic $\ell_2$ datasets, likely because this metric and distribution induce strong metric hubness properties in the dataset and thus a more connected subgraph.

Notably, the Yandex-DEEP benchmark does not exhibit the subgraph connectivity observed in most other benchmarks. Unlike the other embeddings in our collection, Yandex-DEEP was obtained by applying PCA and $\ell_2$ normalization to the last fully-connected layer of GoogleNet [44]. Theoretically, $\ell_2$ normalization should reduce the prevalence of hubs by restricting the distribution to the surface of the unit sphere. Yandex-DEEP also has the most uniform node access distribution of the $\ell_2$ datasets in Figure 9. We believe that this may explain why Yandex-DEEP hubs are not as well-connected.

Table 7 shows the results when hubs are identified using the top-1% threshold. Here, we are able to reject the null hypothesis ($p < 0.05$) in all but one case for both the Mann-Whitney $U$-test and the two-sample $t$-test. We also observe much larger effect sizes with the top-1% selection criterion, demonstrating that the strongest hubs cluster tightly in the graph. Taken together, these two statistical

tests support the idea that hub nodes are often connected to other hubs.

**Table 6: Two-sample $t$-test and Mann-Whitney U-test results. Hub nodes are selected using the P95 threshold of the node access distribution.**

| Dataset | Dim | Mann-Whitney | Two-Sample $t$-Test | Effect Size |
|---|---|---|---|---|
| IID Normal (Angular) | 16 | 03629 | 0.3090 | 0.0267 |
| IID Normal (L2) | 16 | $< 10^{-5}$ | $< 10^{-5}$ | 0.3737 |
| IID Normal (Angular) | 32 | 0.0335 | 0.0516 | 0.0872 |
| IID Normal (L2) | 32 | $< 10^{-5}$ | $< 10^{-5}$ | 0.4275 |
| IID Normal (Angular) | 64 | 0.0216 | 0.0148 | 0.1165 |
| IID Normal (L2) | 64 | $< 10^{-5}$ | $< 10^{-5}$ | 0.3965 |
| IID Normal (Angular) | 128 | 0.0083 | 0.0083 | 0.1284 |
| IID Normal (L2) | 128 | $< 10^{-5}$ | $< 10^{-5}$ | 0.3773 |
| IID Normal (Angular) | 256 | 0.0009 | 0.0007 | 0.1723 |
| IID Normal (L2) | 256 | $< 10^{-5}$ | $< 10^{-5}$ | 0.2620 |
| IID Normal (Angular) | 1024 | 0.1000 | 0.1114 | 0.0652 |
| IID Normal (L2) | 1024 | $< 10^{-5}$ | $< 10^{-5}$ | 0.2361 |
| IID Normal (Angular) | 1536 | 0.0957 | 0.1141 | 0.0645 |
| IID Normal (L2) | 1536 | $< 10^{-5}$ | $< 10^{-5}$ | 0.2512 |
| GloVe | 100 | $< 10^{-5}$ | $< 10^{-5}$ | 0.2550 |
| NYTimes | 256 | $< 10^{-5}$ | $< 10^{-5}$ | 0.4488 |
| GIST | 960 | $< 10^{-5}$ | $< 10^{-5}$ | 0.3645 |
| Yandex-DEEP | 96 | 0.5002 | 0.5000 | 0.0000 |
| Microsoft-SpaceV | 100 | 0.1586 | 0.1585 | 0.0535 |

**Table 7: Two-sample $t$-test and Mann-Whitney U-test results. Hub nodes are selected using the P99 threshold of the node access distribution.**

| Dataset | Dim | Mann-Whitney | Two-Sample $t$-Test | Effect Size |
|---|---|---|---|---|
| IID Normal (Angular) | 16 | 0.0006 | 0.0006 | 0.1745 |
| IID Normal (L2) | 16 | $< 10^{-5}$ | $< 10^{-5}$ | 0.6621 |
| IID Normal (Angular) | 32 | 0.0347 | 0.0347 | 0.0972 |
| IID Normal (L2) | 32 | $< 10^{-5}$ | $< 10^{-5}$ | 0.8173 |
| IID Normal (Angular) | 64 | 0.0359 | 0.0417 | 0.0927 |
| IID Normal (L2) | 64 | $< 10^{-5}$ | $< 10^{-5}$ | 0.8725 |
| IID Normal (Angular) | 128 | 0.0093 | 0.0070 | 0.1316 |
| IID Normal (L2) | 128 | $< 10^{-5}$ | $< 10^{-5}$ | 0.8428 |
| IID Normal (Angular) | 256 | $< 10^{-5}$ | $< 10^{-5}$ | 0.3110 |
| IID Normal (L2) | 256 | $< 10^{-5}$ | $< 10^{-5}$ | 0.8582 |
| IID Normal (Angular) | 1024 | 0.1472 | 0.1318 | 0.0598 |
| IID Normal (L2) | 1024 | $< 10^{-5}$ | $< 10^{-5}$ | 0.8314 |
| IID Normal (Angular) | 1536 | $< 10^{-5}$ | $< 10^{-5}$ | 0.2356 |
| IID Normal (L2) | 1536 | $< 10^{-5}$ | $< 10^{-5}$ | 0.8568 |
| GloVe | 100 | $< 10^{-5}$ | $< 10^{-5}$ | 0.7642 |
| NYTimes | 256 | $< 10^{-5}$ | $< 10^{-5}$ | 0.9305 |
| GIST | 960 | $< 10^{-5}$ | $< 10^{-5}$ | 0.6829 |
| Yandex-DEEP | 96 | 0.0013 | 0.0013 | 0.1614 |
| Microsoft-SpaceV | 100 | 0.0011 | 0.0011 | 0.1644 |

## 6.5 Hub-Highway Nodes Enable Fast Traversal

In this section, we show that not only are highway nodes strongly connected to each other, they also allow queries to quickly traverse the similarity search graph. While it is not surprising that a well-connected subgraph of frequently visited nodes would enable this behavior, it is not immediately clear that queries use the highway in the way predicted by our hypothesis – to quickly identify a neighborhood for deep exploration. To investigate this question, we track the sequence of nodes visited during beam search. We examine thousands of queries to determine the fraction of time spent on hub nodes in different phases of search.

Since beam search takes a variable number of steps for each query, we normalize by the total search length when presenting the results. More formally, suppose that $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_l$ is a length-$l$ sequence of such nodes visited by a query. We use the hub node assignment heuristic discussed in section 6.4 to label $h(\mathbf{x}_i)$ each of these nodes as hubs / non-hubs. Then, we then split the sequence into bins and compute the prevalence of hubs in each bin. For bin $B_i$, this is given by

$$\left(\frac{1}{|B_i|}\right) \sum_{\mathbf{x}_j \in B_i} h(\mathbf{x}_j)$$

where $|B_i|$ is the bin size (fixed to 30 in our analysis). By averaging this value over all queries, we can plot the likelihood of visiting a hub as the search progresses.

Figure 10 shows results for the Gist, GloVe, Microsoft SpaceV and Yandex-DEEP benchmark datasets. We observe that queries tend to concentrate in the highway structures early in search, shown by the high percentage of hub nodes visited in the first 5-10% of the search steps. This result suggests that the highway allows queries to quickly navigate the similarity search graph until they find the region of the graph best suited for deep exploration. The propensity of the query to visit hubs appears to be tied to the hubness properties of the dataset. For example, the GloVe dataset uses the angular distance and has less pronounced hubs (Figure 9), and queries spend a lower percentage of their time in hub nodes in this dataset (Figure 10b). On the other hand, the GIST dataset has some of the highest rates of highway utilization and is also our highest-dimensional $\ell_2$ dataset.



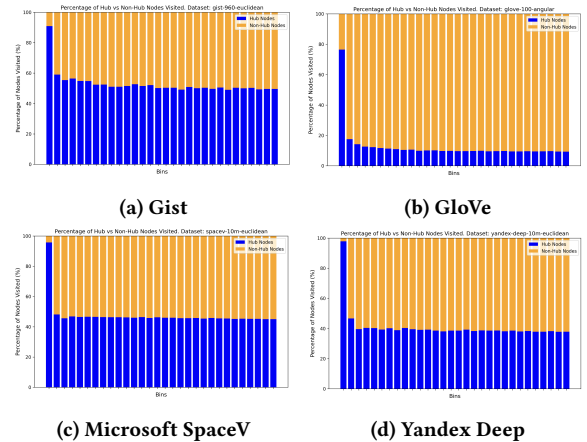| | |
|---|---|
| **(a) Gist** | **(b) GloVe** |
| **(c) Microsoft SpaceV** | **(d) Yandex Deep** |

**Figure 10: Highway nodes allow queries to traverse the graph faster. We observe that early in search,**

## 6.6 Discussion

We observe substantial evidence from the above sequence of experiments supporting the *Hub Highway Hypothesis*. Although it has long been established that the hubness phenomenon negatively affect common applications, such as clustering and even near neighbor search recall, existing ANNS methods, such as HNSW and NSG [51] mostly emphasize algorithmic improvements and performance optimizations only. We have shown that leveraging

the inherent structures in the data, particularly hub-highway occurrences, should be central to the design of new scalable similarity search indexes.

As discussed earlier, the skewness of the $k$-occurrence distribution $N_k$ directly affects the quality of search in terms of recall. Given our empirical results showing that hub-highways allow queries to traverse the graph faster, interesting future research directions might consider leveraging this distribution to improve edge pruning heuristics or judiciously choose which edges to explore during beam search over proximity graphs.

**Small-World Graphs:** The network science research community has known for decades that long-range connections and hubs induce the formation of "small-world" graphs that are easily traversed [48]. This idea has been enormously influential in ANNS, providing the motivation for both NSW and HNSW [30, 31], but our results suggest that ANN graphs constructed over low-dimensional datasets may not in fact exhibit small-world properties. Because the $k$-occurrence distribution is near-uniform for intrinsically low-dimensional data distributions, a pure $k$NN graph (without pruning or long-range links) will not create hubs with a high in-degree. We believe that hierarchical structures are helpful in low dimensions because they help to induce hub behavior, by ensuring that search always begins from a small set of nodes. However, this is not necessary to produce hubs and induce small-world properties in high dimensions. The $k$NN graph construction process is sufficient on its own, because the hub highway emerges in high dimensions.

**Scientific Implications:** Our work reveals that the hub highway is an intrinsic and naturally-forming structure in high-dimensional proximity graphs. We believe this observation is both novel and has important implications for the scaling potential of many traversal heuristics. Specifically, we expect the benefits of sophisticated search initialization methods to decay under hub-style preferential attachment.

Initialization is a recurring and popular research direction for graph-based near-neighbor search, and leading algorithms vary greatly in their initialization techniques. The research question dates back to the seminal 1993 paper by Arya and Mount [2], which conjectured that clever search initialization – in their case, via $kd$-tree – could improve performance over random initialization. Over the following three decades, the research community has investigated diverse stratified sampling based on clusters [41], vantage-point trees [19], seeds formed from the graph expansion of $kd$-trees [21] and previously-visited nodes [47], hierarchical graphs [32], hierarchical clustering [35], dataset medoids [23], and several other methods before finally returning to cluster-stratified candidates [36, 37] and random entry nodes [22].

How is it that these early papers on graph search show performance gains, even as today's best vector databases return to simple initializations without performance loss? Our hub highway hypothesis offers a clear explanation for this apparent contradiction: In the early 2000s and 2010s, datasets were low-dimensional and initialization was important to avoid local minima and long graph detours. However, modern vector databases contain data that is sufficiently high-dimensional to naturally form a fast-routing structure, explaining why initialization no longer drives performance. Based on our results, we conjecture that the largest algorithmic improvements to graph-based ANNS should come from optimizations

that affect the connectivity and cost of traversal in the base graph, such as link pruning and search algorithm design.

## 7 CONCLUSION

Approximate near neighbor search has become an increasingly crucial computational workload in recent years with the seminal Hierarchical Navigable Small World (HNSW) algorithm continuing to garner significant interest and adoption from practitioners. To that end, this paper presents the first comprehensive study on the utility of the hierarchical component of HNSW over numerous large-scale datasets and performance metrics. Ultimately, we find that the hierarchy of HNSW provides no clear benefit on high-dimensional datasets and can be **removed** without any discernible loss in performance while providing memory savings.

While similar observations have been made before in the literature [8, 10, 28], we are, to our knowledge, the first to conduct an exhaustive study over modern benchmark datasets and taking extensive care to compare implementations with performance engineering parity. Furthermore, we go beyond prior works and study *why* the hierarchy does not help, culminating in our introduction of the *Hub Highway Hypothesis*, an empirical result on how proximity graphs built over high-dimensional metric spaces leverage a small subset of well-connected nodes to traverse the network quickly, akin to highways and exits in transportation systems.

As part of our benchmarking study and analysis, we also release an open-source, high-performance flat navigable small world implementation, called `flatnav`, that we hope will serve as an additional resource to both researchers and industry practitioners. We also hope that our work can inspire future directions in algorithm design for ANN search, such as in better strategies for exploiting the hub highway structure we observed or even designing new algorithms that do manage to leverage the idea of hierarchy effectively. In the short term, we believe that our results provide immediate implications for practitioners seeking to save memory or simplify their vector database implementations and we look forward to further partnering with the community on these endeavors.

## REFERENCES

[1] Cecilia Aguerrebere, Ishwar Bhati, Mark Hildebrand, Mariano Tepper, and Theodore Willke. 2023. Similarity Search in the Blink of an Eye with Compressed Indices. *Proceedings of the VLDB Endowment* 16 (08 2023), 3433–3446. https://doi.org/10.14778/3611479.3611537
[2] Sunil Arya and David M Mount. 1993. Approximate nearest neighbor queries in fixed dimensions.. In *SODA*, Vol. 93. Citeseer, 271–280.
[3] Martin Aumüller, Erik Bernhardsson, and Alexander John Faithfull. 2018. ANN-Benchmarks: A Benchmarking Tool for Approximate Nearest Neighbor Algorithms. *Inf. Syst.* 87 (2018). https://api.semanticscholar.org/CorpusID:36089988
[4] Jon Louis Bentley. 1975. Multidimensional binary search trees used for associative searching. *Commun. ACM* 18, 9 (1975), 509–517.
[5] Kevin Beyer, Jonathan Goldstein, Raghu Ramakrishnan, and Uri Shaft. 1999. When is "nearest neighbor" meaningful?. In *Database Theory—ICDT'99: 7th International Conference Jerusalem, Israel, January 10–12, 1999 Proceedings 7*. Springer, 217–235.
[6] Alina Beygelzimer, Sham Kakade, and John Langford. 2006. Cover trees for nearest neighbor. In *Proceedings of the 23rd international conference on Machine learning*. 97–104.
[7] Leonid Boytsov and Bilegsaikhan Naidan. 2013. Engineering efficient and effective non-metric space library. In *Similarity Search and Applications: 6th International Conference, SISAP 2013, A Coruña, Spain, October 2-4, 2013, Proceedings 6*. Springer, 280–293.
[8] Benjamin Coleman, Santiago Segarra, Alexander J Smola, and Anshumali Shrivastava. 2022. Graph reordering for cache-efficient near neighbor search. *Advances in Neural Information Processing Systems* 35 (2022), 38488–38500.

[9] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. 2022. *Introduction to algorithms*. MIT press.

[10] Magdalen Dobson, Zheqi Shen, Guy E Blelloch, Laxman Dhulipala, Yan Gu, Harsha Vardhan Simhadri, and Yihan Sun. 2023. Scaling Graph-Based ANNS Algorithms to Billion-Size Datasets: A Comparative Analysis. *arXiv preprint arXiv:2305.04359* (2023).

[11] Mohamad Dolatshah, Ali Hadian, and Behrouz Minaei-Bidgoli. 2015. Ball*-tree: Efficient spatial indexing for constrained nearest-neighbor search in metric spaces. *arXiv preprint arXiv:1511.00628* (2015).

[12] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library. *arXiv preprint arXiv:2401.08281* (2024).

[13] Rick Durrett. 2019. *Probability: theory and examples*. Vol. 49. Cambridge university press.

[14] Chao Feng, Wuchao Li, Defu Lian, Zheng Liu, and Enhong Chen. 2022. Recommender forest for efficient retrieval. *Advances in Neural Information Processing Systems* 35 (2022), 38912–38924.

[15] Damien François, Vincent Wertz, and Michel Verleysen. 2007. The Concentration of Fractional Distances. *IEEE Transactions on Knowledge and Data Engineering* 19 (2007), 873–886. https://api.semanticscholar.org/CorpusID:13220558

[16] Leo J Guibas and Robert Sedgewick. 1978. A dichromatic framework for balanced trees. In *19th Annual Symposium on Foundations of Computer Science (sfcs 1978)*. IEEE, 8–21.

[17] Ruiqi Guo, Philip Sun, Erik M. Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. 2019. Accelerating Large-Scale Inference with Anisotropic Vector Quantization. In *International Conference on Machine Learning*. https://api.semanticscholar.org/CorpusID:218614141

[18] 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.

[19] Masajiro Iwasaki. 2010. Proximity search in metric spaces using approximate k nearest neighbor graph. *IPSJ Trans. Database* 3, 1 (2010), 18–28.

[20] Masajiro Iwasaki. 2016. Pruned bi-directed k-nearest neighbor graph for proximity search. In *International Conference on Similarity Search and Applications*. Springer, 20–33.

[21] Masajiro Iwasaki and Daisuke Miyazaki. 2018. Optimization of indexing based on k-nearest neighbor graph for proximity search in high-dimensional data. *arXiv preprint arXiv:1810.07355* (2018).

[22] Shikhar Jaiswal, Ravishankar Krishnaswamy, Ankit Garg, Harsha Vardhan Simhadri, and Sheshansh Agrawal. 2022. Ood-diskann: Efficient and scalable graph anns for out-of-distribution queries. *arXiv preprint arXiv:2211.12850* (2022).

[23] 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. *Advances in Neural Information Processing Systems* 32 (2019).

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

[25] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-Scale Similarity Search with GPUs. *IEEE Transactions on Big Data* 7 (2017), 535–547. https://api.semanticscholar.org/CorpusID:926364

[26] Jon M Kleinberg. 2000. Navigation in a small world. *Nature* 406, 6798 (2000), 845–845.

[27] Patrick Lewis, Ethan Perez, Aleksandara Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *ArXiv* abs/2005.11401 (2020). https://api.semanticscholar.org/CorpusID:218869575

[28] Peng-Cheng Lin and Wan-Lei Zhao. 2019. Graph based nearest neighbor search: Promises and failures. *arXiv preprint arXiv:1904.02077* (2019).

[29] Thomas Low, Christian Borgelt, Sebastian Stober, and Andreas Nürnberger. 2013. The hubness phenomenon: Fact or artifact? *Towards Advanced Data Analysis by Combining Soft Computing and Statistics* (2013), 267–278.

[30] Yury Malkov, Alexander Ponomarenko, Andrey Logvinov, and Vladimir Krylov. 2012. Scalable distributed algorithm for approximate nearest neighbor search problem in high dimensional general metric spaces. In *Similarity Search and Applications: 5th International Conference, SISAP 2012, Toronto, ON, Canada, August 9-10, 2012. Proceedings 5*. Springer, 132–147.

[31] 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.

[32] Yury Malkov and Dmitry A. Yashunin. 2016. Efficient and Robust Approximate Nearest Neighbor Search Using Hierarchical Navigable Small World Graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42 (2016), 824–836. https://api.semanticscholar.org/CorpusID:8915893

[33] Henry B Mann and Donald R Whitney. 1947. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics* (1947), 50–60.

[34] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems* 26 (2013).

[35] Javier Vargas Munoz, 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.

[36] Jiongkang Ni, Xiaoliang Xu, Yuxiang Wang, Can Li, Jiajie Yao, Shihai Xiao, and Xuecang Zhang. 2023. Diskann++: Efficient page-based search over isomorphic mapped graph index using query-sensitivity entry vertex. *arXiv preprint arXiv:2310.00402* (2023).

[37] Yutaro Oguri and Yusuke Matsui. 2024. Theoretical and Empirical Analysis of Adaptive Entry Point Selection for Graph-based Approximate Nearest Neighbor Search. *arXiv preprint arXiv:2402.04713* (2024).

[38] William Pugh. 1990. Skip lists: a probabilistic alternative to balanced trees. *Commun. ACM* 33, 6 (1990), 668–676.

[39] Milos Radovanovic, Alexandros Nanopoulos, and Mirjana Ivanovic. 2010. Hubs in space: Popular nearest neighbors in high-dimensional data. *Journal of Machine Learning Research* 11, sept (2010), 2487–2531.

[40] Parikshit Ram and Kaushik Sinha. 2019. Revisiting kd-tree for nearest neighbor search. In *Proceedings of the 25th acm sigkdd international conference on knowledge discovery & data mining*. 1378–1388.

[41] Thomas B Sebastian and Benjamin B Kimia. 2002. Metric-based shape retrieval in large databases. In *2002 International Conference on Pattern Recognition*, Vol. 3. IEEE, 291–296.

[42] Harsha Simhadri, George Williams, Martin Aumüller, Matthijs Douze, Artem Babenko, Dmitry Baranchuk, Qi Chen, Lucas Hosseini, Ravishankar Krishnaswamy, Gopal Srinivasa, Suhas Subramanya, and Jingdong Wang. 2022. Results of the NeurIPS'21 Challenge on Billion-Scale Approximate Nearest Neighbor Search.

[43] Philip Sun, David Simcha, Dave Dopson, Ruiqi Guo, and Sanjiv Kumar. 2024. SOAR: improved indexing for approximate nearest neighbor search. *Advances in Neural Information Processing Systems* 36 (2024).

[44] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, D. Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2014. Going deeper with convolutions. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2014), 1–9. https://api.semanticscholar.org/CorpusID:206592484

[45] Michel Talagrand. 1994. Concentration of measure and isoperimetric inequalities in product spaces. *Publications Mathématiques de l'Institut des Hautes Études Scientifiques* 81 (1994), 73–205. https://api.semanticscholar.org/CorpusID:119668709

[46] Jeffrey Travers and Stanley Milgram. 1977. An experimental study of the small world problem. In *Social networks*. Elsevier, 179–197.

[47] Jingdong Wang and Shipeng Li. 2012. Query-driven iterated neighborhood graph search for large scale indexing. In *Proceedings of the 20th ACM international conference on Multimedia*. 179–188.

[48] Duncan J Watts and Steven H Strogatz. 1998. Collective dynamics of 'small-world' networks. *nature* 393, 6684 (1998), 440–442.

[49] Lihi Zelnik-Manor and Pietro Perona. 2004. Self-Tuning Spectral Clustering. In *Neural Information Processing Systems*. https://api.semanticscholar.org/CorpusID:17066951

[50] Jianshu Zhao, Jean Pierre Both, Luis M Rodriguez-R, and Konstantinos T Konstantinidis. 2024. GSearch: ultra-fast and scalable genome search by combining K-mer hashing with hierarchical navigable small world graphs. *Nucleic Acids Research* 52, 16 (2024), e74–e74.

[51] Xi Zhao, Yao Tian, Kai Huang, Bolong Zheng, and Xiaofang Zhou. 2023. Towards Efficient Index Construction and Approximate Nearest Neighbor Search in High-Dimensional Spaces. *Proc. VLDB Endow.* 16 (2023), 1979–1991. https://api.semanticscholar.org/CorpusID:258718568