

ACORN: Performant and Predicate-Agnostic Search Over Vector Embeddings and Structured Data

Liana Patel
Stanford University
Stanford, USA
lianapat@stanford.edu

Carlos Guestrin
Stanford University
Stanford, USA
guestrin@stanford.edu

Peter Kraft
DBOS, Inc.
USA
peter.kraft@dbos.dev

Matei Zaharia
UC Berkeley
Berkeley, USA
matei@berkeley.edu

ABSTRACT

Applications increasingly leverage mixed-modality data, and must jointly search over *vector data*, such as embedded images, text and video, as well as *structured data*, such as attributes and keywords. Proposed methods for this *hybrid search* setting either suffer from poor performance or support a severely restricted set of search predicates (e.g., only small sets of equality predicates), making them impractical for many applications. To address this, we present ACORN, an approach for performant and predicate-agnostic hybrid search. ACORN builds on Hierarchical Navigable Small Worlds (HNSW), a state-of-the-art graph-based approximate nearest neighbor index, and can be implemented efficiently by extending existing HNSW libraries. ACORN introduces the idea of *predicate subgraph traversal* to emulate a theoretically ideal, but impractical, hybrid search strategy. ACORN’s predicate-agnostic construction algorithm is designed to enable this effective search strategy, while supporting a wide array of predicate sets and query semantics. We systematically evaluate ACORN on both prior benchmark datasets, with simple, low-cardinality predicate sets, and complex multi-modal datasets not supported by prior methods. We show that ACORN achieves state-of-the-art performance on all datasets, outperforming prior methods with 2–1,000× higher throughput at a fixed recall.

CCS CONCEPTS

• **Information systems** → **Information retrieval query processing**; **Data structures**.

KEYWORDS

Vector Search, Approximate Nearest Neighbor Search, Hybrid Search

1 INTRODUCTION

Due to the representation strength of modern deep learning models, vector embeddings have become a powerful first-class datatype for wide-ranging applications that use retrieval-augmented generation [3, 65] or similarity-based search [18, 21, 42]. As a result, vector databases and indices are seeing increasing adoption in many production use cases. These systems provide an efficient approximate-nearest-neighbor (ANN) search interface over embedded unstructured data e.g., images, text, video, or user profiles.

However, many applications must jointly query *both unstructured and structured data*, requiring ANN search in combination with

predicate filtering. For example, customers on an e-commerce site can search for t-shirts similar to a reference image, while filtering on price [64]. Similarly, researchers performing a literature review may search with both natural language queries and filters on publication date, keywords or topics [54]. Likewise, a data scientist working on outlier detection can find misclassified images by retrieving those that look similar to a reference dog but have the label "cat" [2, 7].

To leverage diverse data modalities, applications need data management systems that effectively support *hybrid search queries*, i.e., similarity search with structured predicates. Such systems require (1) **query performance**, i.e., efficient and accurate search despite variance in workload characteristics, such as selectivity, attribute correlations, and scale, and (2) **expressive query semantics**: support for diverse query predicates that may not be known a priori (e.g., user-entered keywords, range searches, or regex matching).

Unfortunately, existing systems fall short of these goals. Three commonly used methods are pre-filtering [62, 64], post-filtering [1, 5, 62, 64, 67], and specialized data structures for low-cardinality predicate sets [25, 49, 63, 66]. *Pre-filtering* first finds all records in the dataset that pass the query predicate, then performs brute force similarity-search over the filtered vector set. This approach scales poorly, becoming inefficient for medium to high selectivity predicates on large datasets. Alternatively, *post-filtering* first searches an ANN index, then removes results that fail the query predicate. Since the database vectors closest to the query vector may not pass the predicate, post-filtering methods must typically expand the search scope. This is often expensive, particularly for search predicates with low selectivity or low correlation to the query vector, as we show in Figure 2. Milvus [62], Weaviate [1], AnalyticDB-V [64], and FAISS-IVF [5] build systems using these two core methods, and suffer from their performance limitations.

Recognizing these limitations, recent works construct *specialized indices* designed for hybrid search workloads with low-cardinality predicate sets consisting of equality predicate operators. For example, Filtered-DiskANN [25] outperforms prior baselines, but restricts the cardinality of the predicate set to about 1,000 and only supports equality predicates. HQANN [66] and NHQ [12] similarly constrain the predicate set to a small number of equality filters and in addition allow only a single structured attributes per dataset entry. These methods are often impractical since many applications have large, or unbounded predicate sets that are unknown a priori. In general, the possible predicate set’s cardinality grows

exponentially with each attribute’s cardinality, which itself may be large. Thus, we instead propose a *predicate-agnostic* index, which can support arbitrary and unbounded predicate sets.

In this paper, we propose **ACORN** (ANN Constraint-Optimized Retrieval Network), a novel approach for performant and predicate-agnostic hybrid search that can serve high-cardinality and unbounded predicate sets. We propose two indices: *ACORN- γ* , designed for high-efficiency search, and *ACORN-1*, designed for low construction overhead in resource-constrained settings. Both methods modify the hierarchical navigable small world (HNSW) index, a state-of-the-art graph-based index for ANN search, and are easy to implement in existing HNSW libraries.

ACORN tackles both the *performance limitations* of pre- and post-filtering, as well as the *semantic limitations* of specialized indices. ACORN proposes the idea of *predicate subgraph traversal* during search. As the name implies, the search strategy traverses the subgraph of the ACORN index induced by the set of nodes that pass the query predicate. ACORN designs the index such that these arbitrary predicate subgraphs approximate an HNSW index. Unlike pre- and post-filtering, this allows ACORN to provide sub-linear retrieval times despite variance in *correlation* between query vectors and predicates, which we find to be a major challenge for existing hybrid search systems. ACORN also serves wide-ranging predicate sets by employing a predicate-agnostic construction that alters HNSW’s algorithm to create a denser graph. Specifically, we introduce a predicate-agnostic neighbor expansion strategy in *ACORN- γ* based on target predicate selectivity thresholds, which can be estimated empirically with or without knowing the predicate set. In conjunction, we propose a predicate-agnostic compression heuristic to efficiently manage the index space footprint while maintaining efficient search. We also explore the trade-off space between search performance and construction overhead, designing *ACORN-1* to approximate *ACORN- γ* ’s search performance while further reducing the time-to-index (TTI) for resource-constrained settings.

We systematically evaluate *ACORN- γ* and *ACORN-1* on four datasets: SIFT1M [35], Paper [63], LAION [55], and TripClick [54]. Our evaluation includes both prior benchmark datasets, with simple, low-cardinality predicate sets, which prior specialized indices can serve, as well as more complex datasets with millions of possible predicates, which existing indices cannot to handle. On each, *ACORN- γ* achieves state-of-the-art hybrid search performance with 2–1000 \times higher queries per second (QPS) at 0.9 recall compared to prior methods. Specifically, ACORN achieves 2-10 \times higher QPS on prior benchmarks, over 30 \times higher QPS on new benchmarks, and over 1,000 \times higher QPS at scale, on a 25-million-vector dataset. We find that *ACORN-1* empirically approximates *ACORN- γ* , attaining at most 5 \times lower QPS at fixed recall but 9–53 \times lower TTI compared to *ACORN- γ* . Our detailed evaluation demonstrates the effectiveness of ACORN’s predicate-subgraph traversal strategy and predicate-agnostic construction techniques.

2 BACKGROUND

Existing methods for Approximate Nearest Neighbor (ANN) search can be broadly categorized as tree-based [15–17, 19, 28, 45, 50, 56], hashing-based [9–11, 24, 26, 29, 30, 40, 41, 44, 46, 52, 59, 69], quantization-based [23, 27, 34, 35, 39], and graph-based [22, 25, 32,

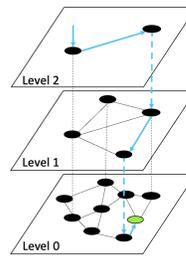


Figure 1: Schematic drawing of search over an HNSW index. The search path is shown by blue arrows, beginning on level 2 and ending on level 0 at the query point, shown in green.

47, 48, 58, 68]. In this work build on HNSW, a graph-based method that is empirically one of the best-performing on high-dimensional datasets, and we adapt it to support hybrid search.

Graph-based ANN methods have gained popularity due their state-of-the-art performance on varied ANN benchmarks [13, 57]. These methods typically perform search using a greedy routing strategy that traverses a graph index, starting from a pre-defined entry point. The index itself forms a proximity graph $G(V, E)$, such that each dataset point is represented by a vertex and edges connect nearby points. The index construction algorithm typically aims to approximate subgraphs of the Delaunay graph [38]. While the Delaunay graph guarantees convergence of a greedy routing algorithm, it is impossible to efficiently construct for arbitrary metric spaces [51]. Thus graph methods focus on more tractable approximations of Delaunay subgraphs, such as the *Relative Neighbor Graph (RNG)* [37, 60], and the *Nearest-Neighbor Graph (NNG)* [8, 20].

2.1 Hierarchical Navigable Small Worlds

As Figure 1 illustrates, HNSW forms a hierarchical, multi-level graph index with bounded degree. Below we briefly summarize the HNSW search and construction algorithm.

The HNSW construction algorithm iteratively inserts each point into the graph index, to construct a navigable graph with bounded degree, specified by parameter M . For each inserted element v , first, a maximum layer index l is stochastically chosen using an exponentially decaying probability distribution, normalized by the constant $m_L = 1/\ln(M)$. The level assignment probability ensures that the expected characteristic path length increases with the level index. Intuitively, the upper-most level contains the longest-range links, which will be traversed first by the search algorithm, and the bottom-most level contains the shortest-range links, which are traversed last by the search algorithm. The insertion procedure then proceeds in two phases. In the first phase, a greedy search is performed iteratively from the top layer, beginning at a pre-defined entrypoint down to the $(l + 1)$ th layer. At each of these levels, the greedy subroutine chooses a single node that becomes the entrypoint into the next layer. In the second phase, the greedy search iterates over level l to level 0. The greedy search at each level now chooses efc nodes as candidate edges. Of these candidates, at most M are selected to become neighbors of v according to an RNG-based pruning algorithm [31]. At level 0, the degree bound is increased $2 \times M$, which is shown to empirically improve performance.

The HNSW search algorithm begins its traversal from a pre-defined entry point at the upper-most layer of the multilayer graph,

Algorithm 1: HNSW-ANN-SEARCH(x_q, K, efs)

Input: query vector x_q , number of nearest neighbors to return K , size of dynamic candidate list efs
Output: K nearest elements to x_q
 $e \leftarrow$ entry-point to hnsw graph
 $W \leftarrow \emptyset$ // set of current nearest
 $L \leftarrow level(e)$ // Top hnsw level
for $l \leftarrow L \dots 1$ **do**
 $e \leftarrow$ SEARCH-LAYER($x_q, e, efs = 1, l$)
end
 $W \leftarrow$ SEARCH-LAYER($x_q, e, efs = efs, l = 0$)
return K nearest elements from W to x_q

illustrated in Figure 1. The traversal then follows an iterative search strategy from the top level downwards. At each level a greedy search is used to choose a single node, which becomes the entry-point into the next level. Once the bottom level is reached, rather than greedily choosing a single node, the search algorithm greedily chooses K nearest elements to return. We outline this process in Algorithm 1. The search parameter efs provides a tradeoff between search quality and efficiency by controlling the size of the dynamic candidate list stored during the bottom level’s greedy search.

3 PROBLEM DEFINITION AND CHALLENGES

In this section we formally define the hybrid search setting and then analyze the performance challenges that existing predicate-agnostic methods, i.e., pre- and post-filtering, face. Our analysis leads us to explore several important workload characteristics. Specifically, we will consider predicate selectivity, the dataset size, and *query correlation*, which we introduce, formally define and find to be a major challenge for post-filtering methods.

We will later leverage our understanding of existing performance challenges in Section 4 to formulate a theoretically ideal hybrid search solution. Then, in Section 7, we will revisit the workload characteristics discussed in this section to rigorously evaluate ACORN’s search performances.

3.1 Hybrid Search Definitions

Let $D = \{e_1, e_2, \dots, e_n\} = \{(x_1, a_1), (x_2, a_2), \dots, (x_n, a_n)\}$ be a dataset consisting of n entities, each with a vector component, $x_i \in \mathbb{R}^d$, and a structured attribute-tuple, a_i , associated with entity e_i . Let $X = \{x_1, x_2, \dots, x_n\}$ denote the set of vectors in the dataset, and $dist(a, b)$ is the metric distance between any two points. Let $A = \{a_1, a_2, \dots, a_n\}$ be the set of structured attributes in the dataset. We will denote $X_p \subseteq X$ as the subset of vectors corresponding to entities in the dataset that pass a given predicate p . We refer to the *selectivity* (s) of predicate p as the fraction of entities from D that satisfy the predicate, where $0 \leq s \leq 1$.

We consider the **hybrid search problem**, described as follows. Given a dataset D , target K , and query $q = (x_q, p_q)$, where $x_q \in \mathbb{R}^d$, and p_q is a predicate, retrieve x_q ’s K nearest neighbors that pass the predicate p_q . We will specifically focus on the problem of *approximate* nearest neighbor search w.r.t x_q . Here, our goal is to maximize both search accuracy and search efficiency. We will measure accuracy by $recall@K = \frac{G \cap R}{K}$, where G is the ground

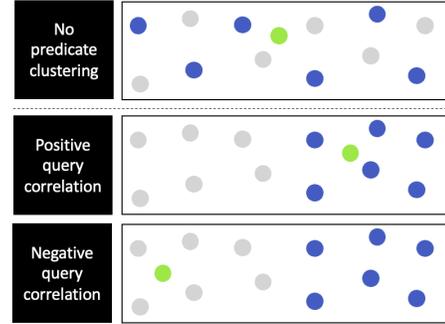


Figure 2: Schematic drawing of a dataset with no predicate clustering (top), a dataset with predicate clustering and positive query correlation (middle), and a dataset with predicate clustering and negative query correlation (bottom). Dark blue circles show points that pass the predicate, and light gray circles show points that fail the predicate. The query vectors are shown in green.

truth set of K nearest neighbors to x_q that satisfy p_q , and R is the retrieved set.

3.2 Search Performance of Baseline Methods

We now analyze the search complexity of two predominant baseline methods, pre- and post-filtering. We will consider how varied workload characteristics impact the search behavior of these methods. Through our analysis, we will make the standard assumption that distance computations dominate search performance. We note that HNSW’s unfiltered search complexity is $O(\log(n) + K)$.

Pre-filtering linearly scans X_p , computing distances over each point that passes the search predicate. This yields a hybrid search complexity of $O(|X_p|) = O(sn + K)$. While pre-filtering always achieves perfect recall, its search complexity scales poorly for large dataset sizes or selectivities, growing linearly in either variable.

Post-filtering, by contrast, performs ANN-search over X to find the closest query vector to x_q , then expands the search scope to find K vectors that pass the query predicate, p . Intuitively, search performance varies depending on *correlation* between the query vector and the vectors in X_p . When the vectors of X_p are close to the query vector, post-filtering over HNSW has a search complexity of $O(\log(n) + K)$. If the vectors in X_p are uniformly distributed within X , then post-filtering’s expected search complexity is $O(\log(n) + K/s)$. However, vectors of X_p may be far away from the query vector, leading to a worst case of $O(n)$ search performance.

We see that the search performance of either baseline is *not robust* to variations in selectivity, dataset size, and query correlation. We empirically verify these limitations in section 7 (figures 9, 10).

3.2.1 Formalizing Query Correlation. We will now formalize the notion of *query correlation*, which we find is key challenge for post-filtering-based systems. As Figure 2 shows, query correlation occurs when the vectors of X_p are non-uniformly distributed in X and instead cluster together relative to the vectors in X . We refer to this phenomenon as *predicate clustering*. When predicate clustering occurs, a query vector may be either close or far away from the predicate cluster containing its search targets, inducing query correlation.

Definition: Query Correlation. We will consider the query-to-target distances for the *given dataset* compared to the expected query-to-target distances for a *hypothetical dataset*, under which no clustering is present. Formally, we define the *query correlation* of the hybrid search workload Q over dataset D as:

$$C(D, Q) = \mathbb{E}_{(x_i, p_i) \in Q} [\mathbb{E}_{R_i} [g(x_i, R_i)] - g(x_i, X_{p_i})]$$

We let R_i be a random set variable of $|X_{p_i}|$ vectors drawn *uniformly* from X , defined for each hybrid query $(x_i, p_i) \in Q$. We define $g(x, S) = \min_{y \in S} \text{dist}(x, y)$ to be the function mapping the query vector x to the minimum distance of neighbors from the given vector set $S \subseteq R^d$. Note that $g(x_i, X_{p_i})$ is the ground-truth hybrid-search target of the query (x_i, p_i) .

If, on average, query vectors are closer to their targets in X_{p_i} , the true dataset of hybrid search targets, than in R_i , the no-clustering dataset, then the workload has *positive query correlation*. If the reverse is true, the workload has *negative query correlation*. We may also consider nearest-neighbor distance rather than the metric distance in the above definition. We also note that we can easily extend this definition to consider K targets of the hybrid search, rather than one, by summing distances over the K search targets.

4 THEORETICAL IDEAL HYBRID SEARCH PERFORMANCE WITH HNSW

For a given hybrid search query, we define the theoretically ideal search performance using HNSW data structures as the performance attainable if we knew the search predicate p_q during construction. In this case, we could construct an HNSW index over X_p . We call this the **oracle partition index** for that query. The complexity of searching this index is $O_s(\log sn + K)$. Notably, the search performance of the oracle partition index outperforms both pre- and post-filtering across variations in predicate selectivity, data size, and query correlation. While pre-filtering’s search scales in $|X_p|$, search over the oracle partition scales *sublinearly* in $|X_p|$. The oracle partition is also robust under variations in query correlation: it does not require the search scope expansion used in post-filtering.

Despite its ideal search performance, the oracle partition index requires us to know all search predicates in advance and to create a full HNSW index per predicate. In practice, the oracle partition index is not possible to construct because query predicate sets are often unknown during construction and have high or unbounded cardinality. Building an HNSW per predicate would require prohibitive amounts of space and time. Thus, in this work, we will instead approximate search over the oracle partition index for a particular query, without ever explicitly constructing this index.

5 ACORN OVERVIEW

We now describe ACORN, a predicate-agnostic approach for state-of-the-art hybrid search. We propose two variants, which we refer to as **ACORN- γ** (5.1, 5.2) and **ACORN-1** (5.3). We design ACORN- γ to achieve efficient search performance, and we design ACORN-1 to approximate ACORN- γ ’s search performance while further reducing the algorithm’s time to index (TTI) and space footprint for resource-constrained settings.

ACORN’s core idea is to search over the index’s *predicate subgraph*, i.e., the subgraph induced by X_p for a given search predicate

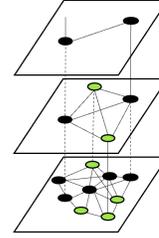


Figure 3: An illustration of the predicate subgraph, shown by the green nodes. ACORN searches over the predicate subgraph to emulate search over an oracle partition index.

Table 1: Summary of Notation

Symbol	Description
γ	neighbor expansion factor for ACORN index
M_β	compression parameter for ACORN index
ef	size of dynamic candidate list in ACORN greedy search
M	degree bound for traversed nodes during ACORN search
$m_L = 1/\ln M$	level normalization constant for ACORN index
e	entry-point to ACORN index
e_p	entry-point to ACORN’s predicate p ’s subgraph
$l(v)$	maximum level index of node v in ACORN index
$N^l(v)$	neighbor list of node v at level l
$N_p^l(v)$	filtered neighbors of node v at level l under predicate p
X_p	vector dataset that passes predicate p
s	selectivity
n	size of dataset

p , as shown in Figure 3. We modify the HNSW construction algorithms so that arbitrary predicate subgraphs emulate an HNSW oracle partition index without the need to explicitly construct one. ACORN- γ achieves this by constructing a denser version of HNSW, which we parameterize by a neighbor list expansion factor, γ , a compression factor, M_β , and the HNSW parameters, efc and M . Then by adding a filter step during search to ignore neighbors that fail the predicate, we find ACORN- γ ’s search can efficiently navigate to and traverse over the predicate subgraph, even under variations in query correlation. Meanwhile, ACORN-1 expands neighbor lists *during search* rather than during construction to *approximate* ACORN- γ ’s dense graph structure without building it.¹

Overall, ACORN prescribes a simple and general framework for performant hybrid search based on the idea of *predicate-subgraph traversal*. The core techniques we propose are *predicate-agnostic* neighbor-list expansions and pruning during construction in combination with predicate-based filtering during search. While this framework can be applied to a variety of graph-based ANN indices, in this work we focus on HNSW due their state-of-the-art performance and widespread use.

5.1 ACORN- γ Search Algorithm

Algorithm 2 outlines the greedy search algorithm ACORN uses at each level, beginning from the top level at a pre-defined entry-point. The main difference between ACORN’s search algorithm and that of HNSW is how neighbor look-ups (line 9) are performed at

¹For highly selective queries where even ACORN’s predicate subgraph would be disconnected within the larger ACORN graph, ACORN falls back to pre-filtering, which is effective for such queries. ACORN is configured with a minimum selectivity, s_{min} , under which it should use pre-filtering when a query is estimated to be more selective than s_{min} . We describe how to configure γ based on s_{min} in Section 5.2.

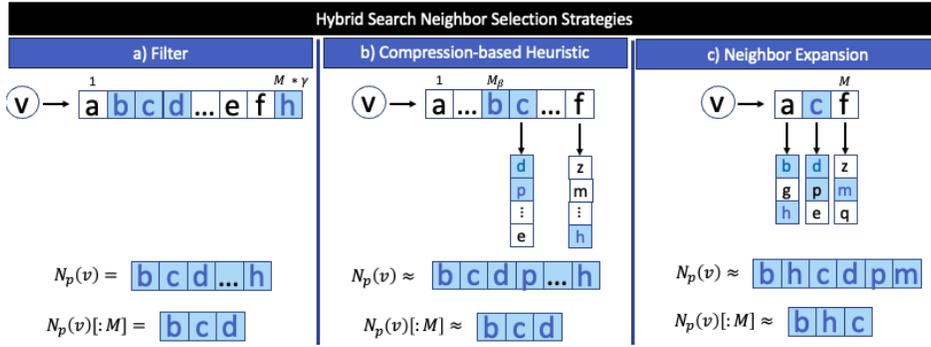


Figure 4: Diagram of ACORN’s neighbor selection strategies. Blue nodes represent neighbors that pass the query predicate. Sub-figure (a) shows the simple predicate-based filter applied to uncompressed edge lists of size $M \cdot \gamma$, followed by truncation to size $M = 3$. Sub-figure (b) shows the compression-based heuristic. Sub-figure (c) shows the neighbor expansion strategy used in ACORN-1.

Algorithm 2: ACORN-SEARCH-LAYER(x_q, p_q, e, γ, l)

Input: query vector x_q , query predicate p_q , entry-point e , number of nearest neighbors to return ef , level to search l
Output: ef nearest elements to x_q

```

1  $T \leftarrow e$  // visited set
2  $C \leftarrow e$  // candidate set
3  $W \leftarrow e$  // dynamic list of found nearest neighbors
4 while  $|C| > 0$  do
5    $c \leftarrow \text{extract } \arg \min_{x \in C} \|x_q - x\|$ 
6    $f \leftarrow \text{get } \arg \max_{x \in W} \|x_q - x\|$ 
7   if  $\text{dist}(c, x_q) > \text{dist}(f, x_q)$  and  $|W| \geq ef$ 
8     break
9    $\text{neighborhood} \leftarrow \text{GET-NEIGHBORS}(c, l, p_q)$ 
10  for each  $v \in \text{neighborhood}[1:M]$ 
11    if  $v \notin T$ 
12       $T \leftarrow T \cup v$ 
13       $f \leftarrow \arg \max_{x \in W} \|x_q - x\|$ 
14      if  $\text{dist}(v, x_q) < \text{dist}(f, x_q)$  or  $|W| < ef$ 
15         $C \leftarrow C \cup v$ 
16         $W \leftarrow W \cup v$ 
17        if  $|W| > ef$ 
18          remove furthest element from  $W$  to  $x_q$ 
19
20 end
21 return  $W$ 

```

each visited node, c . While HNSW simply checks the neighbor list, $N^l(c)$, ACORN performs additional steps to recover an appropriate neighborhood for the given search predicates.

Specifically, ACORN- γ uses two neighbor look-up strategies, a simple filter method, shown in Figure 4(a), and a compression-based heuristic, shown in Figure 4(b), which is compatible with the compression strategy we optionally apply during construction, detailed in Section 5.2. For each visited node, v , the filter-based neighbor look-ups simply scan the neighbor list $N^l(v)$ to find the sub-list of neighbors that pass the predicate, $N_p^l(v)$. If $N_p^l(v)$ contains more than M nodes, we take the first M and return this as v ’s neighborhood. The compression-based neighbor look-ups instead partially expand the neighbor set $N^l(v)$ to include a subset of v ’s two-hop

neighbors, before performing filtering and truncation. This procedure entails two phases. The first phase iterates through the first M_β nodes of $N^l(v)$, simply filtering as in the previous strategy. The second phase iterates over the remainder of the neighbor list, expanding the search neighborhood to include neighbors of neighbors, before again filtering according to the query predicate. M_β is a construction parameter which we will discuss in the next section.

5.2 ACORN- γ Construction Algorithm

We construct the ACORN- γ index by applying two core modifications to the HNSW indexing algorithm: first, we expand each node’s neighbor list, and then we apply a novel predicate-agnostic pruning method to compress the index. Both of these steps are summarized in Figure 5.

Neighbor List Expansion. While HNSW collects M approximate nearest neighbors as candidate edges for each node in the index, ACORN collects $M \cdot \gamma$ approximate nearest neighbors as candidate edges per node. To find these candidates during construction, ACORN uses a metadata-agnostic search over its graph index. Specifically, the neighbor lookup strategy at each node, v , on level l , simply accesses the neighbor list $N^l(v)$ and returns the first M nodes. Note that although each node contains up to $M \cdot \gamma$ neighbors, we assume by construction that M neighbors per node are sufficient for maintaining navigability of the graph index. Thus, considering truncated neighbor lists while traversing the graph allows us to avoid unnecessary distance computations and TTI slowdowns.

One simple choice for γ is $\frac{1}{s_{min}}$, where s_{min} is the minimum predicate selectivity we plan to serve before resorting to pre-filtering. As we discuss in Section 6, ACORN’s indexing time and space footprint increase proportionally to γ . Meanwhile, pre-filtering becomes a competitive baseline at low predicate selectivity values, as we show in Figure 9a. Thus, ACORN is able to balance construction and search efficiency by using pre-filtering as a fall-back for queries with low-selectivity predicates. This leads to a simple *cost-based model* during search: if the estimated predicate selectivity of a given query is greater than $1/\gamma$, search the ACORN- γ index, otherwise pre-filter. We note that leveraging pre-filtering in this way may degrade search efficiency, but not result quality, when errors occur in selectivity estimates. If a query’s true predicate selectivity is above $1/\gamma$, but the estimate is below, the system will mistakenly pre-filter,

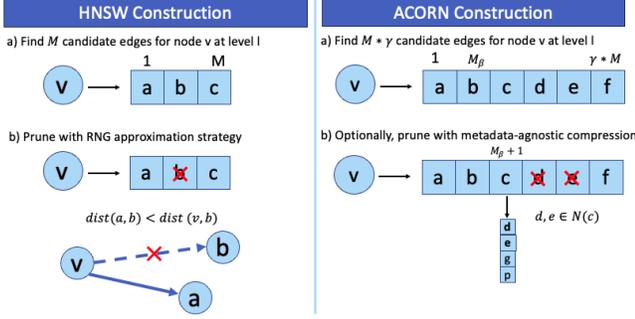


Figure 5: A comparison of HNSW and ACORN- γ 's strategies for (a) selecting candidate edges, shown for $M=3$, and (b) pruning candidate edges for each inserted node v , shown for $M=3$, $M_\beta=2$, $\gamma=2$.

obtaining perfect recall at possibly lower QPS than if the ACORN index was instead searched. If the reverse is true, the system will mistakenly search the ACORN index, whereas pre-filtering would have offered similar QPS and perfect recall.

Compression. A key challenge with ACORN- γ 's neighbor expansion step is that it increases index size and TTI. The increased index size poses a significant issue particularly for memory-resident graph indices, like HNSW. To address this, we introduce a *predicate-agnostic pruning* technique. While we could apply compression to the full index, as discussed in Section 6.1, we specifically target the bottom level's neighbor lists since they contribute most significantly to the indexing overhead. This follows from the exponentially decaying level assignment probability ACORN uses.

The core idea of the pruning procedure is to precisely retain each node's nearby neighbors in the index, while approximating farther away neighbor during search. We use the *tunable compression parameter*, M_β , where $0 \leq M_\beta \leq M \cdot \gamma$. During construction, ACORN chooses each node's final neighbor list by automatically retaining the nearest M_β candidate edges and aggressively pruning the remaining candidates. During search we can recover the first M_β neighbors of each node v directly from the neighbor list $N^l(v)$, and approximate remaining neighbors by looking at 2-hop neighbors during search, as we described in Section 5.1.

Figure 5 outlines this pruning procedure applied to node v 's candidate neighbor list. The algorithm iterates over the ordered candidate edge list and keeps the first M_β candidates. Over the remaining sub-list of candidates, the algorithm applies the following pruning procedure at each node. Let H be the dynamic set of v 's chosen two-hop neighbors, initialized to \emptyset . We prune candidate c if it is contained in H ; otherwise, we keep c and add all of its neighbors to H . The pruning procedure stops after iterating over all candidates, or if $|H|$ plus the number of chosen edges exceeds $M \cdot \gamma$. The pruned and ordered neighbor list is then stored in the ACORN index and H is discarded.

We highlight that the neighbor expansion during search, described in Section 5.1, can recover pruned neighbors regardless of the query predicate. It follows from ACORN's pruning rule that any node x that was pruned from some node v 's neighbor list, $N^l(v)$, must be in the neighbor list $N^l(y)$ such that y is a neighbor of v with index greater than M_β . During search, the neighbor lookup at v on level l will perform a neighbor-list expansion for all neighbors with an index greater than M_β , thus checking $N^l(y)$ and finding x .

We now briefly describe why HNSW's pruning, a *metadata-blind* mechanism, is insufficient for hybrid search. Consider the simple scenario shown in Figure 5. For a node, v , inserted into the HNSW index at an arbitrary level l , the algorithm generates candidates neighbors a , b and c . HNSW's pruning rule iterates over v 's candidate neighbor list in order of nearest to farthest neighbors. Node b is pruned since there exists a neighbor a such that b is closer to a than to v . This RNG-approximation strategy corresponds to pruning the longest edge of the triangle formed by a triplet v, a, b . In this case, we can prune the edge $v - b$ and expect a search path to traverse from v to b via a . The problem with this technique arises when we consider the hybrid search setting for an arbitrary predicate. Say v and b pass a given query predicate p_q , but a does not. Then v, b, a do not form a triangle in the predicate subgraph, and we cannot expect to find the path from v to b through a . As a result, HNSW's pruning mechanism will falsely prune edge $v - b$. If we had complete knowledge of all possible query predicates, we could ensure that we only prune edges of triangles such that all three vertices always exist in the same subset of possible predicate subgraphs. FilteredDiskANN [25] takes this approach by restricting the set of possible query predicates. However, for arbitrary query predicates, ensuring this property holds becomes intractable.

5.3 ACORN-1

We now describe ACORN-1, an alternative approach which aims to approximate ACORN- γ 's search performance, while further minimizing index size and TTI. ACORN-1 achieves this by performing the neighbor expansion step solely during search, rather than during construction, as ACORN- γ does. ACORN-1's construction corresponds to the original HNSW index without pruning. This construction corresponds to ACORN- γ 's construction algorithm, with fixed parameters $\gamma = 1$ and $M_\beta = M$.

ACORN-1's main difference from ACORN- γ during search, is its neighbor lookup strategy. Specifically, at each visited node, v , during greedy search, ACORN-1 uses a full neighbor list expansion to consider all one-hop and two-hop neighbors of v , before applying the predicate filter and truncating the resulting neighbor list to size M . Figure 4(c) outlines this procedure.

6 DISCUSSION

In this section we analyze the ACORN index's space complexity, construction complexity and search performance. We focus our attention on ACORN- γ , since ACORN-1's index construction represents a special case of ACORN- γ for fixed parameters ($\gamma = 1$, $M_\beta = M$), and we empirically show that ACORN-1 search approximates ACORN- γ in Section 7. We note that our analysis in Sections 6.2 and 6.3 considers the complexity scaling of the search procedure under the assumption that we build the exact Delaunay graphs rather than approximate ones.

6.1 Index Size

The average memory consumption per node of the ACORN- γ index is $O(M_\beta + M + m_L \cdot M \cdot \gamma)$, assuming the number of bytes per edge is constant. For comparison, average memory consumption per node for the HNSW index scales $O(M + m_L \cdot M)$. Overall, ACORN- γ increases the bottom-level's memory consumption by $O(M_\beta)$ per

node, and increases the higher levels memory consumption by a factor of γ per node.

To understand ACORN’s memory consumption we evaluate the average number of neighbors stored per node. At level 0, compression is applied to the candidate edge lists of size $M \cdot \gamma$ resulting in neighbor sets of length M_β plus a compressed set which scales $O(M)$. We show this empirically in figure 12. On higher levels, nodes have at most $M \cdot \gamma$ edges. We multiply this by the average number of levels that an element is added to, given by $\mathbb{E}[l + 1] = E[-\ln(\text{unif}(0, 1)) * m_L] = m_L + 1$.

While we specifically target compression to level 0 in this work, because it uses the most space, compression could be applied to more levels in bottom-up order to further reduce the index size for large datasets. Denoting nc as the chosen number of compressed levels, the average memory consumption per node in this generalized case is $O(nc(M_\beta + M) + (m_L - nc)(M \cdot \gamma))$.

6.2 Construction Complexity

For fixed parameters M , M_β and efc , ACORN- γ ’s overall expected construction complexity scales $O(n \cdot \gamma \cdot \log(n) \cdot \log(\gamma))$. Compared to HNSW, which has $O(n \cdot \log(n))$ expected construction complexity, ACORN- γ increases TTI by a factor of $\gamma \cdot \log(\gamma)$ due to the expanded edge lists it generates.

We now describe ACORN’s construction complexity in detail by decomposing it into the following three factors (i) the number of nodes in the dataset, given by n (ii) the expected number of levels searched to insert each node into the index, and (iii) the expected complexity of searching on each level. By design, ACORN’s expected maximum level index scales $O(\log n)$ according to its level-assignment probability, which is the same as HNSW. This provides our bound on (ii).

Turning our attention to (iii), we will first consider the length of the search path and then consider the computation cost incurred at each visited node. For the HNSW level probability assignment, it is known that the expected greedy search path length is bounded by a constant $S = \frac{1}{1 - \exp(-m_L)}$ [48]. We can bound ACORN’s expected search path length by $O(\gamma)$ since the path reaches a greedy minima in a constant number of steps and proceeds to expand the search scope by at most $M \cdot \gamma$ nodes to collect up to $M \cdot \gamma$ candidate neighbors during construction.

The computation complexity at each visited node along the search path is $O(\log(\gamma))$, seen as follows. For each node visited, we first check its neighbor list to find at most M un-visited nodes, on which we perform distance computations in $O(M \cdot d)$ time. Then, we update the sorted lists of candidate nodes and results in $O(M \cdot d \cdot \log(\gamma \cdot M))$ time. Treating M and γ as constants, we see that at each visited node the computation complexity is $O(\log \gamma)$ and for greedy search at each level, the complexity is $O(\gamma \cdot \log(\gamma))$. Multiplying by $n \cdot \log(n)$ yields ACORN’s final expected construction complexity, $O(n \cdot \gamma \cdot \log(n) \cdot \log(\gamma))$.

6.3 Search Analysis

Turning our attention to ACORN- γ ’s search algorithm, we will first point out several properties of HNSW that ACORN’s predicate subgraphs aim to emulate. In Figure 7 we empirically show that ACORN’s search performance approximates that of the HNSW

oracle partition index. We will then describe ACORN’s expected search complexity. We define $l : X \rightarrow \mathbb{N}$ to be the mapping of nodes to their maximum level index in ACORN- γ .

6.3.1 Index and Search Properties. Intuitively, for a given query, ACORN’s predicate subgraph will emulate the HNSW oracle partition index when the predicate subgraph forms a hierarchical structure, each node in the subgraph has degree close to M , the subgraph has a fixed entrypoint at its maximum level index that we can efficiently find during search, and the subgraph is connected. We will examine each of these properties separately and consider when they hold. We also note one main difference between ACORN’s predicate subgraphs and HNSW that arises due to ACORN’s predicate-agnostic pruning: each level of ACORN approximates a KNN graph, while each level of HNSW approximates a RNG graph. While this difference does not affect ACORN’s expected search complexity in Section 6.3.2, Malkov et al. [48] demonstrated that the RNG-based pruning empirically improves performance.

Hierarchy. First, we observe that the arbitrary predicate subgraph $G(X_p)$ forms a controllable hierarchy similar to the HNSW oracle partition index built over X_p with parameter M . This is by design. ACORN- γ ’s construction fixes M , and consequently m_L , the level normalization constant. As a result, nodes of X_p in the ACORN- γ index are sampled at rates equal to the level probabilities of the HNSW partition. Ensuring this level sampling holds allows us to bound the expected greedy search path length at each level by a constant, S , as Malkov et al. [48] previously show.

Bounded Degree. Next, we will describe degree bounds, an important factor that impacts greedy search efficiency and convergence. While HNSW upper bounds the degree of each node by M during construction, ACORN- γ enforces this upper bound during search. This ensures ACORN’s search performs a constant number of distance computations per visited node. We now focus our attention on lower bounding the degree of nodes visited during ACORN- γ ’s search over the predicate subgraph.

If a node in the predicate subgraph has degree much lower than M , this could adversely impact the search convergence and thus recall. For a dataset and query predicate that exhibit no predicate clustering, for any node v in $G(X_p)$,

$$\mathbb{E}[|N_p^l(v)|] = |N^l(v)| \cdot s = \gamma \cdot M \cdot s > M, \forall s > s_{min}$$

This also holds as a lower bounds for datasets with predicate clustering, in which case $Pr(x \in N_p^l(v)) > s, \forall x \in N^l(v)$ where v is a node in the predicate cluster. Thus we will continue our lower bound analysis of node degrees under the worst case assumption of no predicate clustering. Using the binomial concentration inequality with parameter s , and union-bounding over the expected search path length, we show that for the search path $\mathcal{P} = v_1 - \dots - v_y$ over an arbitrary predicate subgraph:

$$Pr\left[\bigcup_{v \in \mathcal{P}} (|N_p(v)| \leq (1 - \delta)M)\right] \leq O(\log n \cdot \exp(-\delta^2 \gamma M s / 2))$$

We also analyze the probability that the subgraph traversal gets disconnected, which we bound by:

$$Pr\left[\bigcup_{v \in \mathcal{P}} (|N_p(v)| \leq 0)\right] \leq O(\log n \cdot (1 - s)^{M \cdot \gamma})$$

We see that both bounds decay exponentially in γ .

Fixed Entry-point. Similar to HNSW, ACORN’s search begins from a fixed entry-point, chosen during construction. This pre-defined entry-point provides a simple and effective strategy that is also predicate-independent and robust to variations in *query correlation*, as we empirically show in Figure 10.

Intuitively, we expect the search to successfully navigate from ACORN’s fixed entry-point, e , to the predicate-subgraph entry-point, e_p , when we find a node that passes the predicate on an upper level of the index that is fully connected. In this case, there will exist a one-hop path from e to e_p . We consider e_p to be an arbitrary node that passes a given predicate p and is on the maximum level of the predicate subgraph. The index’s neighbor expansion parameter, γ , causes the index’s upper levels to be denser and, specifically those with less than $M \cdot \gamma$ nodes, to be fully connected. When these fully connected levels contain at least one node that passes the predicate, the search is guaranteed to route from e to e_p . Since ACORN samples all nodes with equal probability at each level, the probability that nodes passing a given predicate, p , exist on some level is proportional to the predicate’s selectivity, which takes a lower bound of $s_{min} = 1/\gamma$.

Connectivity. We note that neither HNSW nor ACORN provides theoretical guarantees on connectivity over its level graphs for arbitrary datasets. Thus we instead rely primarily on empirical results for our analysis. However, for some cases, we can expect ACORN’s predicate subgraph to be connected when the HNSW oracle partition is connected. Two such cases are when X_p exhibits no predicate clustering, or X_p is clustered around a single region. In either case, each node has an expected degree of at least M and each level approximates a KNN graph, which is connected when $K \gg \log n$. We empirically show in Figure 13a that ACORN’s predicate subgraphs exhibit connectivity for real datasets and hybrid search queries. To analyze potential connectivity problems, we recommend benchmarking ACORN’s hybrid search performance against HNSW’s ANN search performance using equivalent M and efc parameters. If a significant gap in accuracy exists, we recommend incrementally increasing γ from its initial value of $1/s_{min}$.

6.3.2 Search Complexity. ACORN- γ ’s expected search complexity scales:

$$O((d + \gamma) \cdot \log(s \cdot n) + \log(1/s))$$

This approximates the HNSW oracle partition’s expected search complexity, $O(d \cdot \log(s \cdot n))$. Intuitively, ACORN- γ ’s search path performs some filtering at the upper levels before likely entering and traversing the predicate sub-graph, during which ACORN incurs a small overhead compared to HNSW search in order to perform the predicate filtering step over each neighbor list.

We derive ACORN- γ ’s search complexity by considering two stages of its search traversal. In the first stage, search begins from a pre-defined entry-point e , which need not pass the query predicate. In this stage, the search performs filtering only, dropping down each level on which the filtered neighbor list, $N_p(e)$, is found to be empty. Once the traversal reaches the first node, e_p that passes the predicate, it enters the second stage, beginning its traversal over the predicate subgraph $G(X_p)$.

In stage 1 the greedy search path on each layer has length 1, and occurs over $O(\log n - \log(s \cdot n))$ expected levels, yielding the

complexity $O(\log(1/s))$. We see this because the expected maximum level index of the full ACORN index graph scales $O(\log n)$ based on its level-assignment probability [48]. Meanwhile, the predicate subgraph $G(X_p)$ of size $s \cdot n$ has an expected maximum level index of $O(\log(s \cdot n))$, once again according to its level sampling procedure.

The second stage of the search traverses the predicate subgraph in expected $O((d + \gamma) \cdot \log(s \cdot n))$ complexity. As we previously describe, the expected maximum level index of the predicate subgraph scales $O(\log(s \cdot n))$. At each level, the expected greedy path length can be bounded by a constant S due to the index level sampling procedure employed during construction. For each node visited along the greedy path, we perform distance computations in $O(d)$ time on at most M neighbors, and perform a constant-time predicate evaluations over at most $M \cdot \gamma$ neighbors.

7 EVALUATION

We evaluate ACORN through a series of experiments on real and synthetic datasets. Overall, our results show the following:

- ACORN- γ achieves state-of-the-art hybrid search performance, outperforming existing methods by 2-1,000× higher QPS at 0.9 recall on both prior benchmark datasets with simple, low-cardinality predicate sets, and more complex datasets with high-cardinality predicate sets. Specifically, ACORN achieves 2-10× higher QPS on prior benchmarks, over 30× higher QPS on new benchmarks, and over 1,000× higher QPS at scale on a 25-million-vector dataset.
- ACORN- γ and ACORN-1 are predicate-agnostic methods, providing robust search performance under variations in predicate operators, predicate selectivity, query correlation, and dataset size.
- ACORN-1 and ACORN- γ exhibit trade-offs between search performance and construction overhead. While ACORN- γ achieves up to 5× higher QPS than ACORN-1 at fixed recalls, ACORN-1 can be constructed with 9-53× lower time-to-index (TTI).

We now discuss our results in detail. We first describe the datasets (7.1) and baselines (7.2) we use. Then, we present a systematic evaluation of ACORN’s search performance (7.3). Finally, we assess ACORN’s construction efficiency (7.4). We run all experiments on an AWS m5d.24xlarge instance with 370 GB of RAM, 96 vCPUs, and 196 threads.

7.1 Datasets

We conduct our experiments on two datasets with low-cardinality predicate sets (LCPS) and two datasets with high-cardinality predicate sets (HCPS). The LCPS datasets allow us to benchmark prior works that only support a constrained set of query predicates. The HCPS datasets consist of more complex and realistic query workloads, allowing us to more rigorously evaluate ACORN’s search performance. Table 2 provides a concise summary of all datasets.

³On the TripClick dataset, we create two distinct query workloads, described in Section 7.1.2. The average selectivity for either workload is .17 (keywords), and .26 (dates).

³On the LAION dataset, we create four distinct query workloads, described in Section 7.1.2. These workloads have average selectivities of .10 (no-cor), .13 (pos-cor), .069 (neg-cor), .056 (regex).

Table 2: Datasets

	Base Data				Query Workload		
	# Vectors	Vector Dim	Vector Source Data	Structured Data	Predicate Operators	Avg. Query Selectivity	Predicate Cardinality
SIFT1M	1,000,000	128	images	random int.	$\text{equals}(y)$	0.083	12
Paper	2,029,997	200	passages	random int.	$\text{equals}(y)$	0.083	12
TripClick	1,055,976	768	passages	clinical area list & publication date	$\text{contains}(y_1 \vee y_2 \vee \dots) \&$ $\text{between}(y_1, y_2)$	0.17, 0.36 ²	$> 10^8$
LAION (1M)	1,000,448	512	images	text captions & keyword list	$\text{regex-match}(y) \&$ $\text{contains}(y_1 \vee y_2 \vee \dots)$	0.056 - 0.13 ³	$> 10^{11}$
LAION (25M)	24,653,427	512	same as above	same as above	same as above	same as above	same as above

7.1.1 Datasets with Low Cardinality Predicate Sets. We use SIFT1M [35] and Paper [63], the two largest publically-available datasets used to evaluate recent specialized indices [25, 63]. For both datasets, we follow related works [25, 62, 63] to generate structured attributes and query predicates: for each base vector, we assign a random integer in the range 1 – 12 to represent structured attributes; and for each query vector, the associated query predicate performs an exact match with a randomly chosen integer in the attribute value domain. The resulting query predicate set has a cardinality of 12.

SIFT1M: The SIFT1M dataset was introduced by Jegou et al. in 2011 for ANN search. It consists of a collection of 1M base vectors, and 10K query vectors. All of the vectors are 128-dimensional local SIFT descriptors [43] from INRIA Holidays images [33].

Paper: Introduced by Wang et al. in 2022, the Paper dataset consists of about 2M base vectors and 10K query vectors. The dataset is generated by extracting and embedding the textual content from an in-house corpus of academic papers.

7.1.2 Datasets with High Cardinality Predicate Sets. We use the TripClick and LAION datasets in our experiments with HCPS datasets.

TripClick: The TripClick dataset, introduced by Rekabsaz et al. in 2021 for text retrieval, contains a *real* hybrid search query workload and base dataset from the click logs of a health web search engine. Each query consists of natural language search terms along with optional filters on clinical areas (e.g. "cardiology", "infectious disease", "surgery") and publication years. Each entity in the base dataset consists of a text passage, with a list of associated clinical areas and a publication date. The dataset contains 28 unique clinical areas and publication dates ranging from 1900 to 2020, resulting in over 2^{28} possible query predicates total. We construct two query workloads, one consisting of queries that used date filters (dates) and another consisting of queries that used clinical area filters (areas). We generate 768-dimensional vectors from the query texts and passage texts using DPR [36], a widely-used, pre-trained encoder for open-domain Q&A. The resulting dataset has about 1M base vectors, and we use a random sample of 1K queries for each query workload.

LAION: The LAION dataset [55] consists of 400M image embeddings and captions describing each image. The vector embeddings are generated from web-scraped images using CLIP [53], a multi-modal language-vision model. In our evaluation, we construct two base datasets using 1M and 25M LAION subsets, both consisting of image vectors and text captions as a structured attribute. We also generate an additional structured attribute consisting of a keyword

list. We assign each image embedding its keyword list by taking the 3 words with highest text-to-image CLIP scores from a candidate list of 30 common adjectives and nouns (e.g., "animal", "scary").

To evaluate a series of micro-benchmarks, we generate four query workloads. For each query workload, we sample 1K vectors from the dataset as query vectors. We construct the regex query workload with predicates that perform regex-matching over the image captions. For each query predicate, we randomly choose strings of 2-10 regex tokens (e.g., " $^{\wedge}[\theta-9]$ "). In addition, we construct three query workloads with predicates, similar to TripClick, that take a keyword list and filter out entities that do not have at least one matching keyword. Using this setup, we are able to easily control for correlation in the workload, and we generate a no correlation (no-cor), positive correlation (pos-cor), and negative correlation (neg-cor) workload. Figure 6 demonstrates some example queries and multi-modal retrieval results taken from each.

7.2 Benchmarked Methods

We briefly overview the methods we benchmark along with tested parameters. We implement ACORN- γ , ACORN-1, pre-filtering, and HNSW post-filtering in C++ in the FAISS codebase [5].

HNSW Post-filtering: To implement HNSW post-filtering, for each hybrid query with predicate selectivity s , we over-search the HNSW index, gathering K/s candidate results before applying the query filter. We note that this differs to some prior work [25], where HNSW post-filtering is implemented by collecting only K candidate results, leading to significantly worse baseline query performance than ours. For the SIFT1M, Paper and LAION datasets, we use the FAISS's default HNSW construction parameters: $M = 32$, $efc = 40$. For the TripClick dataset, we find that the HNSW index for these parameters is unable to obtain high recalls for the standard ANN search task, thus we perform parameter tuning, as is standard. We perform a grid search for $M \in \{32, 64, 128\}$ and $efc \in \{40, 80, 120, 160, 200\}$ and choose the pair that obtains the highest QPS at 0.9 Recall for ANN search. For TripClick, we choose $M = 128$, $efc = 200$. We generate each recall-QPS curve by varying the search parameter efs from 10 to 800 in step sizes of 50.

Pre-filtering: We implement pre-filtering by first generating a list of dataset entries that pass the query predicate and then performing brute force search using FAISS's optimized implementation for distance comparisons. We also efficiently implement all contains predicate evaluations using bitsets since the corresponding structured attributes have low cardinality.

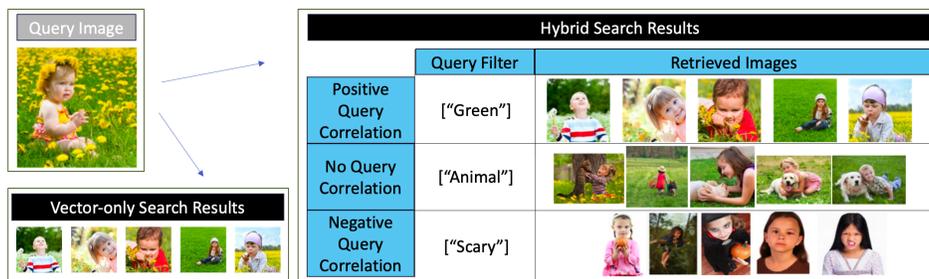


Figure 6: The figure contrasts retrieval results using vector-only similarity search (bottom left) versus hybrid search (right) on the LAION dataset. Both use the same query image (top left), and the hybrid search queries also include a structured query filter consisting of a keyword list, here containing a single keyword. The table on the right shows examples from three hybrid search query workloads: positive query correlation (top), no query correlation (middle), and negative query correlation (bottom).

Filtered-DiskANN: We evaluate both algorithms implemented in FilteredDiskANN [4], namely FilteredVamana and StitchedVamana. For both, we follow the recommended construction and search parameters according to the hyper-parameter tuning procedure described by Gollapudi et al. [25]. For FilteredVamana, we use construction parameters $L = 90, R = 96$, which generated the Pareto-Optimal recall-QPS curve from a parameter sweep over $R \in \{32, 64, 96\}$ and L between 50 and 100. For StitchedVamana, we use construction parameters $R_{small} = 32, L_{small} = 100, R_{stitched} = 64$ and $\alpha = 1.2$, which generated the Pareto-Optimal recall-QPS curve from a parameter sweep over $R_{small}, R_{stitched} \in \{32, 64, 96\}$ and L_{small} between 50 and 100. To generate the recall-QPS curves we vary L from 10 to 650 in increments of 20 for FilteredVamana, and L_{small} from 10 to 330 in increments of 20 for StitchedVamana.

NHQ: We evaluate the two algorithms, NHQ-NPG_NSW and NHQ-NPG_KGraph, proposed in [63]. For both we use the recommended parameters in the released codebase [12]. These parameters were selected using a hyperparameter grid search in order to generate the Pareto-optimal recall-QPS curve for either algorithm on the SIFT1M and Paper datasets. We generate the recall-QPS curve by varying L between 10 and 310 in steps of 20. In Figures 8b and 7b, we show the query performance of KGraph, the more performant of the two algorithms.

Milvus: We test four Milvus algorithms: IVF-Flat, IVF-SQ8, HNSW, and IVF-PQ [6]. For each we test the same parameters as Gollapudi et al. [25]. Since we find that the four Milvus algorithms achieve similar search performance, for simplicity, Figures 8b and 7b show only the method with Pareto-Optimal recall-QPS performance.

Oracle Partition Index: We implement this method by constructing an HNSW index for each possible query predicate in the LCPS datasets. For a given hybrid query, we search the HNSW partition corresponding to the query’s predicate. To construct each HNSW partition and generate the recall-QPS curve, we use the same parameters as the HNSW post-filtering method, described above.

ACORN- γ : We choose the construction parameters M and efc to be the same as the HNSW post-filtering baseline, described above. We find that ACORN- γ ’s search performance is relatively insensitive to the choice of the construction parameter M_β , as Figure 12c shows. Thus, to maintain modest construction overhead, we choose M_β to be a small multiple of M , i.e., $M_\beta = M$ or $M_\beta = 2M$, picking the parameter for each dataset that obtains higher QPS at 0.9 Recall. Specifically, we constrain the memory budget of the index

to be no larger than the Vamana indices on the LCPS datasets and no larger than twice the size of the flat indices for HCPS datasets. We use M_β values of 32 for LAION-1M and LAION-25M, 64 for SIFT1M, Paper, and 128 for TripClick. We choose the construction parameter γ according to the expected minimum selectivity query predicates of each dataset i.e., $\gamma = 12$ for SIFT1M and Paper, $\gamma = 30$ for LAION, and $\gamma = 80$ for TripClick. To generate the recall-QPS curve, we follow the same procedure described above for HNSW post-filtering.

ACORN-1: We construct ACORN-1 and generate the recall-QPS curve following the same procedure we use for ACORN- γ , except that we fix $\gamma = 1$ and $M_\beta = M$.

7.3 Search Performance Results

We will begin our evaluation with benchmarks on the LCPS datasets, on which we are able to run all baseline methods as well as the oracle partition method. We will then present an evaluation on the HCPS datasets. On these datasets, the FilteredDiskANN and NHQ algorithms fail because they assume are unable to handle the high cardinality query predicate sets and non-equality predicate operators. As of this writing, we also find that Milvus cannot support regex-match predicates and contains predicates over variable length lists. As a result, we instead focus on comparing ACORN to the pre- and post-filtering baselines for the HCPS datasets. We report QPS averaged over 50 trials.

7.3.1 Benchmarks on LCPS Datasets. Figure 7 shows that ACORN- γ achieves state-of-the-art hybrid search performance and best approximates the theoretically ideal oracle partition strategy on the SIFT1M and Paper datasets. Notably, even compared to NHQ and FilteredDiskANN, which specialize for LCPS datasets, ACORN- γ consistently achieves 2-10 \times higher QPS at fixed recall values, while maintaining generality. Additionally, we see ACORN-1 approximates ACORN- γ ’s search performance, attaining about 1.5-5 \times lower QPS than ACORN- γ across a range of recall values.

To further investigate the relative search efficiency of ACORN- γ and ACORN-1, we turn our attention to Table 3, which shows the number of distance computations required of either method to obtain Recall@10 equal to 0.8. We see that the oracle partition method is the most efficient, requiring the fewest number of distance computations on both datasets. ACORN- γ is the next most efficient according to number of distance computations. While

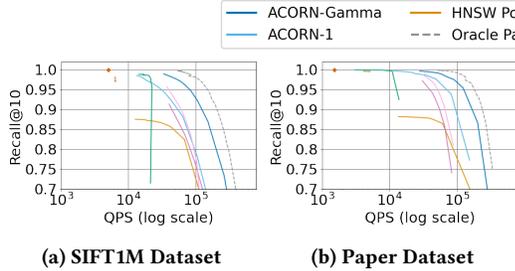


Figure 7: Recall@10 vs QPS on SIFT1M and Paper

ACORN- γ approximates the oracle partition method, its predicate-agnostic design precludes the same RNG-based pruning used to construct the oracle partitions. Rather than approximating RNG-graphs, ACORN- γ 's levels approximate KNN-graphs, which are less efficient to search over explaining the performance gap. The table additionally shows that ACORN-1 is less efficient than ACORN- γ , which is explained by the candidate edge generation used in ACORN-1. While the ACORN- γ index stores up to $M \times \gamma$ edges per node during construction, ACORN-1 stores only up to M edges per node during construction, and *approximates* an edge list of size $M * \gamma$ for each node during search using its neighbor expansion strategy. This approximation results in slight degradation to neighbor list quality and thus search performance. Finally, we see from the table, that HNSW post-filtering is the least efficient of the listed methods. This is because while ACORN-1 and ACORN- γ almost exclusively traverse over nodes that pass the query predicates, the post-filtering algorithm is less discriminating and wastes distance computations on nodes failing the query predicate.

Returning to Figure 7, we see that the relative search efficiency, measured by QPS versus recall, of the oracle partition method, ACORN- γ , and ACORN-1 is not only affected by distance computations, but is also affected by vector dimensionality. We see that both ACORN-1 and ACORN- γ perform closer to the oracle partition method on the Paper dataset, while the performance gap grows slightly on SIFT1M. This is due to the cost of performing the filtering step over neighbor lists during search, which, relative to the cost of distance computations, is higher on SIFT1M than Paper since SIFT1M uses slightly lower-dimensional vectors.

7.3.2 Benchmarks on HCPS Datasets. Figure 8 shows that ACORN outperforms the baselines by 30 – 50 \times higher QPS at 0.9 recall on TripClick and LAION-1M, and as before, ACORN-1 approximates ACORN- γ 's search performance. On both datasets, pre-filtering is prohibitively expensive, obtaining perfect recall at the price of efficiency. Meanwhile, post-filtering fails to obtain high recall, likely

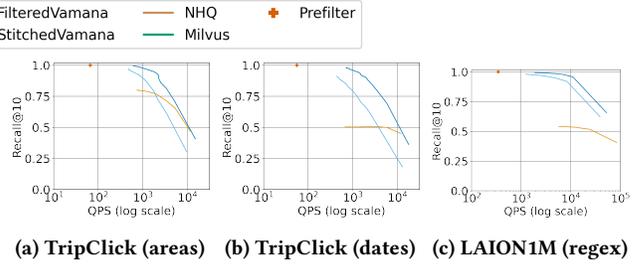


Figure 8: Recall@10 vs QPS on TripClick and LAION-1M

due to the presence of varied query correlation and predicate selectivity, which we further explore further next.

Varied Predicate Selectivity: We use the Tripclick dataset to evaluate ACORN's search performance across a range of realistic predicate selectivities. Figure 9 demonstrates that for each predicate selectivity percentile, ACORN- γ achieves 5-50 \times higher QPS at 0.9 recall compared to the next best-performing baseline. Once again ACORN-1 trails behind ACORN- γ . We see that for low selectivity predicates, the pre-filtering method is most competitive, while the post-filtering baseline suffers from over 10 \times lower QPS than ACORN at fixed recall. However, for high selectivity predicates, pre-filtering becomes less competitive while the post-filtering baseline obtains higher throughput, although its recall remains low.

Varied Query Correlation: Next we control for query correlation and evaluate ACORN on three different query workloads using the LAION-1M dataset. Figure 10 demonstrates that ACORN- γ is robust to variations in query correlation and attains 28-100 \times higher QPS at 0.9 recall than the next best baseline in each case. In the negative correlation case, the performance gap between post-filtering and the ACORN methods is the largest since post-filtering cannot successfully route towards nodes that pass the predicate. In the positive correlation case, ACORN- γ once again outperforms the baselines, but post-filtering become more competitive, although it is still unable to attain recall above 0.9. The pre-filtering method's QPS remains relatively unchanged, and is only affected by small variations in predicate selectivity for each query workload. As before, ACORN-1 approaches ACORN- γ 's search performance.

Scaling Dataset Size: Figure 11 shows ACORN's search performance on LAION-25M with the no-correlation query workload, demonstrating that the performance gap between ACORN and existing baselines only grows as the dataset size scales. At 0.9 recall, ACORN- γ achieves over three orders of magnitude higher QPS than the next best-performing baseline. As before, ACORN-1's search performance approximates that of ACORN- γ .

Table 3: # Distance Computations to Achieve 0.8 Recall

	SIFT 1M	Paper
Oracle Partition	398.0	281.1
ACORN- γ	611.0 (+53.5%)	383.7 (+36.6%)
ACORN-1	999.6 (+151.0%)	567.8 (+101.2%)
HNSW Post-filter	1837.8 (+362.6%)	1425.5 (+406.2%)

* Percentage difference is shown in parenthesis and is relative to oracle partition method

7.4 Index Construction

We will now evaluate ACORN's construction procedure, including its indexing time and space footprint, ACORN- γ 's compression procedure, and the predicate subgraph quality resulting from ACORN- γ 's neighbor expansion approach.

7.4.1 TTI and Space Footprint. First, we analyze ACORN's space footprint and indexing time. Table 4 and 5 show the time-to-index and index size of ACORN- γ and ACORN-1 compared to the best-performing baselines. The reported index sizes for each method show the total space footprint of both vector storage and the index

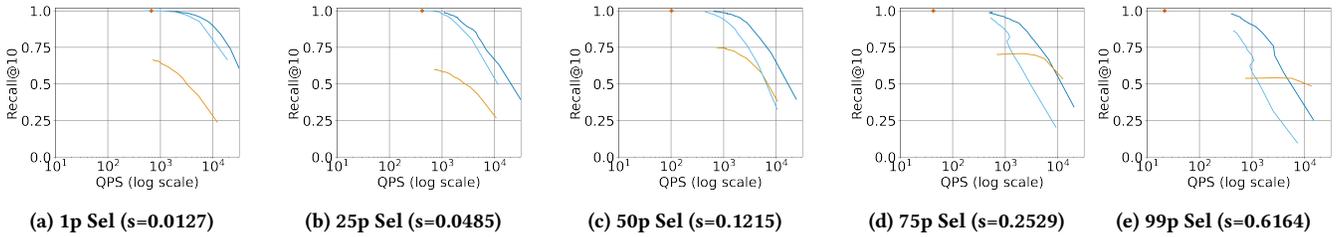


Figure 9: Recall@10 vs QPS for Varied Selectivity Query Filters on TripClick

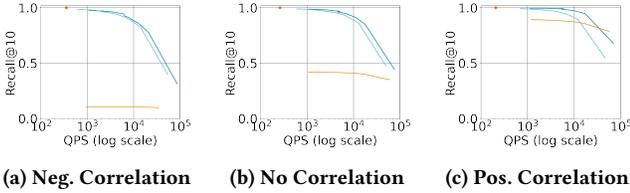


Figure 10: Recall@10 vs QPS on LAION1M

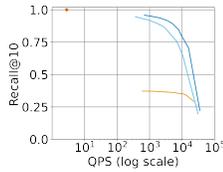


Figure 11: Recall@10 vs QPS on LAION-25M

itself. All methods are measured using the parameters reported in Section 7.2.

We first consider ACORN- γ 's construction overhead. Table 4 shows that across all datasets, ACORN- γ 's TTI is at most $11\times$ higher than HNSW's, and at most $2.15\times$ higher than that of StitchedVamana, the best performing specialized index. Table 5 shows that ACORN- γ 's index size is at most $1.3\times$ larger than that of HNSW, and at least 25% smaller than that of StitchedVamana. The reason for ACORN- γ 's increased index size and TTI compared to HNSW is its candidate-edge generation step during construction, which expands each neighbor list. Meanwhile, ACORN-1 achieves the lowest TTI of all listed baselines in table 4, and its index size is at most $1.25\times$ HNSW's index size and at least 25% smaller than StitchedVamana's index size. We see that while ACORN- γ achieves superior search performance by leveraging a neighbor-list expansion during *construction*, ACORN-1 provides a close approximation at lower TTI and space footprint by instead performing the neighbor-list expansion during *search*. The two algorithms exhibit a trade-off between search performance and construction overhead.

7.4.2 ACORN- γ Pruning. Given ACORN- γ 's higher construction overhead, we investigate the efficiency of its predicate-agnostic compression strategy in reducing index construction costs while maintaining search performance. First, Table 6 shows ACORN- γ 's average out-degree per level for each dataset, confirming that compression on level 0 leads to significantly smaller neighbor lists, compared to level without compression, which may have neighbor lists as large as $M \cdot \gamma$.

Turning our attention to Figure 12, we evaluate three different pruning strategies applied to ACORN- γ 's neighbor lists during

Table 4: TTI (s)

	TripClick	LAION-1M	LAION-25M	Sift1M	Paper
ACORN- γ	9902.9	835.8	38,007.5	148.9	255.6
ACORN-1	322.9	25.9	705.3	8.6	27.0
HNSW	891.0	32.9	1,147.2	11.3	29.2
FilteredVamana	NA	NA	NA	18.3	51.9
StitchedVamana	NA	NA	NA	69.2	189.7

Table 5: Index Size (GB)

	TripClick	LAION-1M	LAION-25M	Sift1M	Paper
ACORN- γ	4.9	2.4	59	0.98	2.5
ACORN-1	4.6	2.3	59	0.93	2.4
HNSW	4.1	2.2	54	.75	2.1
Flat Index	3.1	1.9	47	.51	1.6
FilteredVamana	NA	NA	NA	.61	1.8
StitchedVamana	NA	NA	NA	1.3	3.5

Table 6: ACORN- γ Average Out Degree

	TripClick	LAION-1M	LAION-25M	Sift1M	Paper
Level 0 (compressed)	191	50.1	49.4	87.5	86.0
Level 1	8,075	960	960	384	384
Level 2	54.0	919	937	363	359
Level 3	0	25.3	689	25.3	57.4
Level 4	NA	0	16	0	1.0
$M \cdot \gamma$	10,240	960	960	384	384
M_β	128	32	32	64	64

construction: **i)** ACORN's predicate-agnostic pruning strategy at varied levels of compression indicated by different M_β (Mb) values, where $Mb = 768$ represents no pruning, and lower values represent more aggressive pruning, **ii)** a metadata-aware RNG-based pruning approach, which is employed by FilteredDiskANN's algorithms, and **iii)** HNSW's metadata-blind pruning. We consider TTI, space footprint, the number of candidate edges pruned per node and search performance. The figure represents space footprint measured by the average out degree of nodes on level 0, the level on which each pruning strategy is applied. In addition, the figure shows search performance measured by recall at 20,000 QPS. We note that the recall ranges of the recall-QPS curve generated by different pruning methods varied significantly, leading us to choose a QPS threshold rather than a recall threshold. Interestingly, Figure 12 shows that ACORN's pruning can significantly reduce both the TTI and space footprint by aggressively pruning candidate edges, while maintaining search performance. In comparison, applying HNSW pruning to the index results in significantly degraded hybrid search performance. Meanwhile the metadata-aware RNG-base pruning results in similar search performance to ACORN- γ 's pruning, but it is less efficient by TTI and space footprint than ACORN's pruning for small values of M_β (e.g., $M_\beta = 32, 64$).

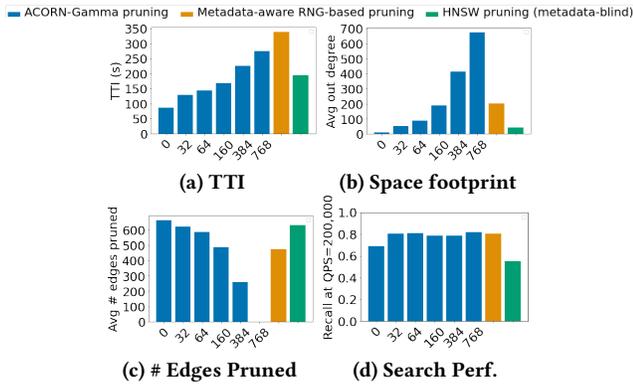


Figure 12: Comparison of pruning methods on SIFT1M and their impact on TTI (a), space footprint of the index (b), the number of candidate edges pruned (c) and search performance (d). M_β values used for ACORN- γ are shown along the x-axis.

7.4.3 Graph Quality. Finally, we investigate the graph quality of ACORN- γ 's predicate subgraphs. Figure 13 compares graph connectivity, graph height, and out degrees for HNSW oracle partitions and ACORN- γ predicate subgraphs across varied predicate selectivities on the TripClick dataset's real hybrid search queries.

From Figure 13a, we see ACORN- γ 's predicate-subgraph connectivity empirically matches or exceeds that of the HNSW oracle partition across selectivities, demonstrating the effectiveness of ACORN- γ 's neighbor expansion strategy. Next, Figure 13b shows that the controlled hierarchy of ACORN- γ 's predicate subgraphs emulate that of the HNSW oracle partitions. Malkov et al. show that HNSW search performance is sensitive to graph height [48]; thus, this result helps explain ACORN- γ 's ability to emulate the search efficiency of the oracle partition. Lastly, Figure 13c examines the average out degree resulting from performing the search-time filtering, described in Figure 4(a), over the ACORN- γ index. We note that sufficiently high, but bounded, out-degrees are important for emulating HNSW's navigability properties, as discussed in Section 6.3. The figure confirms that ACORN's predicate subgraphs have average out-degrees consistently close to and bounded by M . As expected, the HNSW oracle partition has significantly lower average out-degrees than nodes on ACORN- γ 's uncompressed levels because HNSW applies RNG-based pruning. We also note, that the ACORN predicate subgraph with 1 percentile selectivity has lower average out degrees than the other predicate subgraphs because the low selectivity predicates result in fewer than 128 nodes on the largest uncompressed levels, thus capping the maximum out degree per node below $M = 128$. Overall, we observe that ACORN- γ produces high quality predicate subgraphs, which empirically emulate several HNSW properties related to search efficiency.

8 RELATED WORK

Pre- & Post-filtering-based Systems. Many hybrid search systems rely on pre- and post-filtering. While several systems have developed pre-processing methods to perform *faster filtering* during search, these systems fail to reduce the *excessive and expensive distance computations* which bottleneck performance. Weaviate [1] creates an inverted index for structured data ahead of time, then uses it at query time to create a bitmap of eligible candidates during

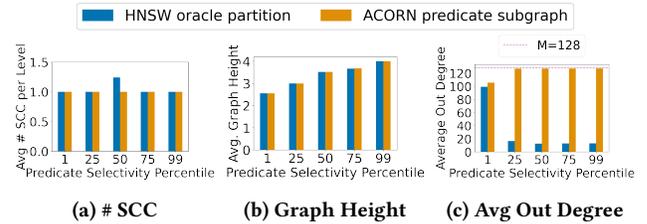


Figure 13: Graph quality of ACORN- γ predicate subgraph evaluated by (a) average number of strongly connected components per level, (b) graph height, and (c) average out degree of nodes across uncompressed levels. Results are shown for the TripClick dataset with 1, 25, 50, 75, and 99 percentile selectivity predicates to generate the predicate subgraph and HNSW oracle partition.

post-filtering. Milvus [62] likewise creates an approved list of points by maintaining a distribution of attributes over the dataset in order to map commonly used query filters to a list of approved points before performing pre- or post-filtering. Several space-partitioning indices like FAISS-IVF [14, 34] and LSH [10] store metadata information in the index, allowing them to rapidly filter entities during post-filtering. Despite the optimized filtering steps in each of these approaches, the core problems of pre- and post-filtering remain, particularly for low correlation or selectivity predicates.

Specialized Indices. Alternatively, several recent works develop novel graph-based algorithms for hybrid search, often improving performance for a constrained set of predicates. NHQ [63] encodes attributes alongside vectors, and then uses a "fusion distance" during search that accounts for vector distances as well as attribute matches. This approach supports only equality query predicates and assumes each dataset entity has only one structured attribute. Filtered-DiskANN [25] proposes two algorithms: FilteredVamana and StitchedVamana. Both methods constrain the query filter cardinality to about 1,000 with only equality predicates so that the index construction steps can use this knowledge to appropriately generate and prune candidate edge lists. Similarly HQI [49] optimizes batch query-processing by assuming a limited cardinality of 20 query predicates to design an efficient partitioning scheme. On the other hand, Qdrant [61] proposes to densify an HNSW graph and perform a filtered greedy search. While this approach aligns intuitively with ACORN's neighbor-list expansions during construction, Qdrant's proposal inadvertently flattens the graph by directly increasing the HNSW parameter M , which impacts HNSW's level normalization constant. Malkov et al. show that HNSW's performance is sensitive to its number of levels, and flattening the graph degrades search performance [48]. In addition, Qdrant's proposed method does not provide a solution for dealing with the increased memory overhead after creating a denser HNSW.

9 CONCLUSION

We proposed ACORN, the first approach for efficient hybrid search across vectors and structured data that supports large and diverse sets of query predicates. ACORN uses a simple, yet effective, search strategy based on the core idea of *predicate subgraph traversal*. We presented two indices, ACORN- γ and ACORN-1, that implement this search strategy by modifying the HNSW indexing algorithm. Our results show that ACORN achieves state-of-the-art hybrid search performance on both prior benchmarks, involving simple,

low-cardinality query predicate sets, as well as more complex benchmarks involving new predicate operators and high cardinality predicate sets. Across both types of benchmarks, ACORN- γ achieves 2–1,000 \times higher QPS at 0.9 recall than prior methods, and ACORN-1 approximates ACORN- γ 's search performance with 9–53 \times lower TTI for resource-constrained settings.

ACKNOWLEDGMENTS

The authors would like to thank Peter Bailis for his valuable feedback on this work.

This research was supported in part by affiliate members and other supporters of the Stanford DAWN project, including Meta, Google, and VMware, as well as Cisco, SAP, and a Sloan Fellowship. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the sponsors.

REFERENCES

- [1] [n. d.]. Filtered Vector Search | Weaviate - vector database. <https://weaviate.io/developers/weaviate/concepts/prefiltering>
- [2] [n. d.]. Pre-label and enrich data with bulk classifications. <https://labelbox.ghost.io/blog/pre-label-and-enrich-your-data-with-bulk-classifications/>
- [3] [n. d.]. Q&A over Documents - LlamaIndex 0.8.43. <https://gpt-index.readthedocs.io/en/latest/>
- [4] 2023. DiskANN. <https://github.com/microsoft/DiskANN> original-date: 2020-06-18T06:18:06Z.
- [5] 2023. Faiss. <https://github.com/facebookresearch/faiss>
- [6] 2023. Milvus Documentation. <https://github.com/milvus-io/milvus-docs> original-date: 2020-05-27T09:12:23Z.
- [7] 2023. visual-layer/fastdup. <https://github.com/visual-layer/fastdup>
- [8] Ann Arbor Algorithms. 2023. KGraph: A Library for Approximate Nearest Neighbor Search. <https://github.com/aaalgo/kgraph> original-date: 2015-05-29T12:38:24Z.
- [9] Alexandr Andoni and Piotr Indyk. 2008. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. *Commun. ACM* 51, 1 (Jan. 2008), 117–122. <https://doi.org/10.1145/1327452.1327494>
- [10] Alexandr Andoni, Piotr Indyk, Thijs Laarhoven, Ilya Razenshteyn, and Ludwig Schmidt. 2015. Practical and optimal LSH for angular distance. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1 (NIPS'15)*. MIT Press, Cambridge, MA, USA, 1225–1233.
- [11] Alexandr Andoni and Ilya Razenshteyn. 2015. Optimal Data-Dependent Hashing for Approximate Near Neighbors. In *Proceedings of the forty-seventh annual ACM symposium on Theory of Computing (STOC '15)*. Association for Computing Machinery, New York, NY, USA, 793–801. <https://doi.org/10.1145/2746539.2746553>
- [12] AshenOn3. 2023. NHQ: An Efficient and Robust Framework for Approximate Nearest Neighbor Search with Attribute Constraint. <https://github.com/AshenOn3/NHQ> original-date: 2021-09-09T08:28:21Z.
- [13] Martin Aumüller, Erik Bernhardsson, and Alexander Faithfull. 2020. ANN-Benchmarks: A benchmarking tool for approximate nearest neighbor algorithms. *Information Systems* 87 (Jan. 2020), 101374. <https://doi.org/10.1016/j.is.2019.02.006>
- [14] Dmitry Baranchuk, Artem Babenko, and Yury Malkov. 2018. Revisiting the Inverted Indices for Billion-Scale Approximate Nearest Neighbors. <https://doi.org/10.48550/arXiv.1802.02422> arXiv:1802.02422 [cs].
- [15] Jon Louis Bentley. 1975. Multidimensional binary search trees used for associative searching. *Commun. ACM* 18, 9 (Sept. 1975), 509–517. <https://doi.org/10.1145/361002.361007>
- [16] Erik Bernhardsson. [n. d.]. annoy: Approximate Nearest Neighbors in C++/Python optimized for memory usage and loading/saving to disk. <https://github.com/spotify/annoy>
- [17] Alina Beygelzimer, Sham Kakade, and John Langford. 2006. Cover trees for nearest neighbor. In *Proceedings of the 23rd international conference on Machine learning (ICML '06)*. Association for Computing Machinery, New York, NY, USA, 97–104. <https://doi.org/10.1145/1143844.1143857>
- [18] Fedor Borisjuk, Siddarth Malreddy, Jun Mei, Yiqun Liu, Xiaoyi Liu, Piyush Maheshwari, Anthony Bell, and Kaushik Rangadurai. 2021. VisRel: Media Search at Scale. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining (KDD '21)*. Association for Computing Machinery, New York, NY, USA, 2584–2592. <https://doi.org/10.1145/3447548.3467081>
- [19] Sanjoy Dasgupta and Yoav Freund. 2008. Random projection trees and low dimensional manifolds. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*. ACM, Victoria British Columbia Canada, 537–546. <https://doi.org/10.1145/1374376.1374452>
- [20] Wei Dong, Charikar Moses, and Kai Li. 2011. Efficient k-nearest neighbor graph construction for generic similarity measures. In *Proceedings of the 20th international conference on World wide web (WWW '11)*. Association for Computing Machinery, New York, NY, USA, 577–586. <https://doi.org/10.1145/1963405.1963487>
- [21] Ming Du, Arnau Ramisa, Amit Kumar K C, Sampath Chanda, Mengjiao Wang, Neelakandan Rajesh, Shasha Li, Yingchuan Hu, Tao Zhou, Nagashri Lakshminarayana, Son Tran, and Doug Gray. 2022. Amazon Shop the Look: A Visual Search System for Fashion and Home. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22)*. Association for Computing Machinery, New York, NY, USA, 2822–2830. <https://doi.org/10.1145/3534678.3539071>
- [22] Cong Fu, Chao Xiang, Changxu Wang, and Deng Cai. 2019. Fast approximate nearest neighbor search with the navigating spreading-out graph. *Proceedings of the VLDB Endowment* 12, 5 (Jan. 2019), 461–474. <https://doi.org/10.14778/3303753.3303754>
- [23] Tiezheng Ge, Kaiming He, Qifa Ke, and Jian Sun. 2014. Optimized Product Quantization. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36, 4 (April 2014), 744–755. <https://doi.org/10.1109/TPAMI.2013.240> Conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [24] Aristides Gionis, Piotr Indyk, and Rajeev Motwani. 1999. Similarity Search in High Dimensions via Hashing. In *Proceedings of the 25th International Conference on Very Large Data Bases (VLDB '99)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 518–529.
- [25] Siddharth Gollapudi, Neel Karia, Varun Sivashankar, Ravishankar Krishnaswamy, Nikit Begwani, Swapnil Raz, Yiyong Lin, Yin Zhang, Neelam Mahapatro, Premkumar Srinivasan, Amit Singh, and Harsha Vardhan Simhadri. 2023. Filtered-DiskANN: Graph Algorithms for Approximate Nearest Neighbor Search with Filters. In *Proceedings of the ACM Web Conference 2023*. ACM, Austin TX USA, 3406–3416. <https://doi.org/10.1145/3543507.3583552>
- [26] Long Gong, Huayi Wang, Mitsunori Ogihara, and Jun Xu. 2020. iDEC: indexable distance estimating codes for approximate nearest neighbor search. *Proceedings of the VLDB Endowment* 13, 9 (May 2020), 1483–1497. <https://doi.org/10.14778/3397230.3397243>
- [27] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. 2020. Accelerating large-scale inference with anisotropic vector quantization. In *Proceedings of the 37th International Conference on Machine Learning (ICML '20, Vol. 119)*. JMLR.org, 3887–3896.
- [28] Michael E. Houle and Michael Nett. 2015. Rank-Based Similarity Search: Reducing the Dimensional Dependence. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37, 1 (Jan. 2015), 136–150. <https://doi.org/10.1109/TPAMI.2014.2343223> Conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [29] 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 (STOC '98)*. Association for Computing Machinery, New York, NY, USA, 604–613. <https://doi.org/10.1145/276698.276876>
- [30] Omid Jafari, Parth Nagarkar, and Jonathan Montaña. 2020. mMLSH: A Practical and Efficient Technique for Processing Approximate Nearest Neighbor Queries on Multimedia Data. In *Similarity Search and Applications (Lecture Notes in Computer Science)*, Shin'ichi Satoh, Lucia Vadicamo, Arthur Zimek, Fabio Carrara, Ilaria Bartolini, Martin Aumüller, Björn Þór Jónsson, and Rasmus Pagh (Eds.). Springer International Publishing, Cham, 47–61. https://doi.org/10.1007/978-3-030-60936-8_4
- [31] J.W. Jaromeczyk and G.T. Toussaint. 1992. Relative neighborhood graphs and their relatives. *Proc. IEEE* 80, 9 (Sept. 1992), 1502–1517. <https://doi.org/10.1109/5.163414> Conference Name: Proceedings of the IEEE.
- [32] Suhas Jayaram Subramanya, Fnu Devvrit, Harsha Vardhan Simhadri, Ravishankar Krishnaswamy, and Rohan Kadekodi. 2019. DiskANN: Fast Accurate Billion-point Nearest Neighbor Search on a Single Node. In *Advances in Neural Information Processing Systems*, Vol. 32. Curran Associates, Inc. https://papers.nips.cc/paper_files/paper/2019/hash/09853c7fb1d3f8e67a61b6bf4a7f8e6-Abstract.html
- [33] Herve Jegou, Matthijs Douze, and Cordelia Schmid. 2008. Hamming Embedding and Weak Geometric Consistency for Large Scale Image Search. In *Computer Vision - ECCV 2008 (Lecture Notes in Computer Science)*, David Forsyth, Philip Torr, and Andrew Zisserman (Eds.). Springer, Berlin, Heidelberg, 304–317. https://doi.org/10.1007/978-3-540-88682-2_24
- [34] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-scale similarity search with GPUs. <http://arxiv.org/abs/1702.08734> arXiv:1702.08734 [cs].
- [35] Hervé Jégou, Matthijs Douze, and Cordelia Schmid. 2011. Product Quantization for Nearest Neighbor Search. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33, 1 (Jan. 2011), 117–128. <https://doi.org/10.1109/TPAMI.2010.57> Conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [36] Vladimir Karpukhin, Barlas Öguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. <https://arxiv.org/abs/2004.04906v3>

- [37] Philip M. Lankford. 1969. Regionalization: Theory and Alternative Algorithms. *Geographical Analysis* 1, 2 (1969), 196–212. <https://doi.org/10.1111/j.1538-4632.1969.tb00615.x> eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1538-4632.1969.tb00615.x>.
- [38] D. T. Lee and B. J. Schachter. 1980. Two algorithms for constructing a Delaunay triangulation. *International Journal of Computer & Information Sciences* 9, 3 (June 1980), 219–242. <https://doi.org/10.1007/BF00977785>
- [39] V. Lempitsky and A. Babenko. 2012. The inverted multi-index. *IEEE Computer Society*, 3069–3076. <https://doi.org/10.1109/CVPR.2012.6248038> ISSN: 1063-6919.
- [40] Mingjie Li, Ying Zhang, Yifang Sun, Wei Wang, Ivor W. Tsang, and Xuemin Lin. 2020. I/O Efficient Approximate Nearest Neighbour Search based on Learned Functions. *2020 IEEE 36th International Conference on Data Engineering (ICDE)* (April 2020), 289–300. <https://doi.org/10.1109/ICDE48307.2020.00032> Conference Name: 2020 IEEE 36th International Conference on Data Engineering (ICDE) ISBN: 9781728129037 Place: Dallas, TX, USA Publisher: IEEE.
- [41] Wanqi Liu, Hanchen Wang, Ying Zhang, Wei Wang, Lu Qin, and Xuemin Lin. 2021. EI-LSH: An early-termination driven I/O efficient incremental c-approximate nearest neighbor search. *The VLDB Journal* 30, 2 (March 2021), 215–235. <https://doi.org/10.1007/s00778-020-00635-4>
- [42] Yiding Liu, Weixue Lu, Suqi Cheng, Daiting Shi, Shuaiqiang Wang, Zhicong Cheng, and Dawei Yin. 2021. Pre-trained Language Model for Web-scale Retrieval in Baidu Search. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining (KDD '21)*. Association for Computing Machinery, New York, NY, USA, 3365–3375. <https://doi.org/10.1145/3447548.3467149>
- [43] David G. Lowe. 2004. Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision* 60, 2 (Nov. 2004), 91–110. <https://doi.org/10.1023/B:VISI.0000029664.99615.94>
- [44] Kejing Lu and Mineichi Kudo. 2020. R2LSH: A Nearest Neighbor Search Scheme Based on Two-dimensional Projected Spaces. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*. 1045–1056. <https://doi.org/10.1109/ICDE48307.2020.00095> ISSN: 2375-026X.
- [45] Kejing Lu, Hongya Wang, Wei Wang, and Mineichi Kudo. 2020. VHP: approximate nearest neighbor search via virtual hypersphere partitioning. *Proceedings of the VLDB Endowment* 13, 9 (May 2020), 1443–1455. <https://doi.org/10.14778/3397230.3397240>
- [46] Qin Lv, William Josephson, Zhe Wang, Moses Charikar, and Kai Li. 2017. Intelligent probing for locality sensitive hashing: multi-probe LSH and beyond. *Proceedings of the VLDB Endowment* 10, 12 (Aug. 2017), 2021–2024. <https://doi.org/10.14778/3137765.3137836>
- [47] Yury Malkov, Alexander Ponomarenko, Andrey Logvinov, and Vladimir Krylov. 2014. Approximate nearest neighbor algorithm based on navigable small world graphs. *Information Systems* 45 (Sept. 2014), 61–68. <https://doi.org/10.1016/j.is.2013.10.006>
- [48] Yu A. Malkov and D. A. Yashunin. 2018. Efficient and robust approximate nearest neighbor search using Hierarchical Navigable Small World graphs. <http://arxiv.org/abs/1603.09320> arXiv:1603.09320 [cs].
- [49] Jason Mohoney, Anil Pacaci, Shihabur Rahman Chowdhury, Ali Mousavi, Ihab F. Ilyas, Umar Farooq Minhas, Jeffrey Pound, and Theodoros Rekatsinas. 2023. High-Throughput Vector Similarity Search in Knowledge Graphs. <http://arxiv.org/abs/2304.01926> arXiv:2304.01926 [cs].
- [50] Marius Muja and David G. Lowe. 2014. Scalable Nearest Neighbor Algorithms for High Dimensional Data. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36, 11 (Nov. 2014), 2227–2240. <https://doi.org/10.1109/TPAMI.2014.2321376> Conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [51] Gonzalo Navarro. 2002. Searching in metric spaces by spatial approximation. *The VLDB Journal* 11, 1 (Aug. 2002), 28–46. <https://doi.org/10.1007/s007780200060>
- [52] Yongjoo Park, Michael Cafarella, and Barzan Mozafari. 2015. Neighbor-sensitive hashing. *Proceedings of the VLDB Endowment* 9, 3 (Nov. 2015), 144–155. <https://doi.org/10.14778/2850583.2850589>
- [53] 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. <https://doi.org/10.48550/arXiv.2103.00020> arXiv:2103.00020 [cs].
- [54] Navid Rekasaz, Oleg Lesota, Markus Schedl, Jon Brassey, and Carsten Eickhoff. 2021. TripClick: The Log Files of a Large Health Web Search Engine. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2507–2513. <https://doi.org/10.1145/3404835.3463242> arXiv:2103.07901 [cs].
- [55] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarezyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. LAION-400M: Open Dataset of CLIP-Filtered 400 Million Image-Text Pairs. <https://doi.org/10.48550/arXiv.2111.02114> arXiv:2111.02114 [cs].
- [56] Chanop Silpa-Anan and Richard Hartley. 2008. Optimised KD-trees for fast image descriptor matching. *IEEE Computer Society*, 1–8. <https://doi.org/10.1109/CVPR.2008.4587638>
- [57] Harsha Vardhan Simhadri, George Williams, Martin Aumüller, Matthijs Douze, Artem Babenko, Dmitry Baranchuk, Qi Chen, Lucas Hosseini, Ravishankar Krishnaswamy, Gopal Srinivasa, Suhas Jayaram Subramanya, and Jingdong Wang. 2022. Results of the NeurIPS'21 Challenge on Billion-Scale Approximate Nearest Neighbor Search. <http://arxiv.org/abs/2205.03763> arXiv:2205.03763 [cs].
- [58] Aditi Singh, Suhas Jayaram Subramanya, Ravishankar Krishnaswamy, and Harsha Vardhan Simhadri. 2021. FreshDiskANN: A Fast and Accurate Graph-Based ANN Index for Streaming Similarity Search. <https://doi.org/10.48550/arXiv.2105.09613> [cs].
- [59] Narayanan Sundaram, Aizana Turmukhmetova, Nadathur Satish, Todd Mostak, Piotr Indyk, Samuel Madden, and Pradeep Dube. 2013. Streaming similarity search over one billion tweets using parallel locality-sensitive hashing. *Proceedings of the VLDB Endowment* 6, 14 (Sept. 2013), 1930–1941. <https://doi.org/10.14778/2556549.2556574>
- [60] Godfried T. Toussaint. 1980. The relative neighbourhood graph of a finite planar set. *Pattern Recognition* 12, 4 (Jan. 1980), 261–268. [https://doi.org/10.1016/0031-3203\(80\)90066-7](https://doi.org/10.1016/0031-3203(80)90066-7)
- [61] Andrei Vasnetsov. [n. d.]. Filtrable HNSW - Qdrant. <https://qdrant.tech/articles/filtrable-hnsw/>
- [62] Jianguo Wang, Xiaomeng Yi, Rentong Guo, Hai Jin, Peng Xu, Shengjun Li, Xi-angyu 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 *Proceedings of the 2021 International Conference on Management of Data (SIGMOD '21)*. Association for Computing Machinery, New York, NY, USA, 2614–2627. <https://doi.org/10.1145/3448016.3457550>
- [63] Mengzhao Wang, Lingwei Lv, Xiaoliang Xu, Yuxiang Wang, Qiang Yue, and Jiongkang Ni. 2022. Navigable Proximity Graph-Driven Native Hybrid Queries with Structured and Unstructured Constraints. <http://arxiv.org/abs/2203.13601> arXiv:2203.13601 [cs].
- [64] Chuangxian Wei, Bin Wu, Sheng Wang, Renjie Lou, Chaoqun Zhan, Feifei Li, and Yuanzhe Cai. 2020. AnalyticDB-V: a hybrid analytical engine towards query fusion for structured and unstructured data. *Proceedings of the VLDB Endowment* 13, 12 (Aug. 2020), 3152–3165. <https://doi.org/10.14778/3415478.3415541>
- [65] Brie Wolfson. 2023. Building chat langchain. <https://blog.langchain.dev/building-chat-langchain-2/>
- [66] Wei Wu, Junlin He, Yu Qiao, Guoheng Fu, Li Liu, and Jin Yu. 2022. HQANN: Efficient and Robust Similarity Search for Hybrid Queries with Structured and Unstructured Constraints. <http://arxiv.org/abs/2207.07940> arXiv:2207.07940 [cs].
- [67] Qianxi Zhang, Shuotao Xu, Qi Chen, Guoxin Sui, Jiadong Xie, Zhizhen Cai, Yaoqi Chen, Yinxuan He, Yuqing Yang, Fan Yang, Mao Yang, and Lidong Zhou. 2023. {VBASE}: Unifying Online Vector Similarity Search and Relational Queries via Relaxed Monotonicity. 377–395. <https://www.usenix.org/conference/osdi23/presentation/zhang-qianxi>
- [68] Weijie Zhao, Shulong Tan, and Ping Li. 2020. SONG: Approximate Nearest Neighbor Search on GPU. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*. 1033–1044. <https://doi.org/10.1109/ICDE48307.2020.00094> ISSN: 2375-026X.
- [69] Bolong Zheng, Xi Zhao, Lianggui Weng, Nguyen Quoc Viet Hung, Hang Liu, and Christian S. Jensen. 2020. PM-LSH: A fast and accurate LSH framework for high-dimensional approximate NN search. *Proceedings of the VLDB Endowment* 13, 5 (Jan. 2020), 643–655. <https://doi.org/10.14778/3377369.3377374>