

# Stratified Graph Indexing for Efficient Search in Deep Descriptor Databases

M M Mahabubur Rahman

toufik@txstate.edu

Texas State University

San Marcos, Texas, USA

Jelena Tešić

jtesic@txstate.edu

Texas State University

San Marcos, Texas, USA

## ABSTRACT

Searching for unseen objects in extensive visual archives is challenging, demanding efficient indexing methods that can support meaningful similarity retrievals. This research paper presents the Stratified Graph (SG) approach for indexing similar deep descriptors by sorting them into distance-sensitive layers based. The indexing algorithm incrementally constructs a bi-directional  $m$ -nearest neighbor graph within each layer, with additional 1-nearest neighbor links from outer layers, providing a distant scaling property in the graph structure. The search process starts from the innermost layer, and the same layer neighbors contribute to enhancing recall, while the distant scaling property enhances search speed, maintaining logarithmic complexity scaling. We compare and contrast SG with six state-of-the-art retrieval methods in four deep-descriptor and two classical-descriptor databases, and we show that the Stratified Graph (SG) indexing and search has smaller memory usage (up to four times), and farther precision and recall improves up to 8% than state-of-art for all six datasets at five retrieval depths.

## KEYWORDS

Similarity search, high-dimensional indexing, deep descriptors, information retrieval

### ACM Reference Format:

M M Mahabubur Rahman and Jelena Tešić. 2023. Stratified Graph Indexing for Efficient Search in Deep Descriptor Databases. In *Proceedings of ACM Multimedia (ACMMM '23)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/XXXXXX.XXXXXXX>

## 1 INTRODUCTION

The world we live in today is a treasure trove of video archives that houses an immense amount of valuable information. To tap into the true potential of these archives, it's crucial to identify objects that are similar to one another as object labeling is sparse and far apart. Achieving this task requires a similarity search in deep descriptors databases that can detect comparable objects that are widely spread throughout the dataset. In this paper we propose to answer the task of *finding the appearances of the new unlabeled object in the existing video archives?*, formulate the task at hand as as the *k-Nearest Neighbor task in the deep descriptor space*.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

ACMMM '23, October 29 - November 3, 2023, Ottawa, Canada

© 2023 Association for Computing Machinery.

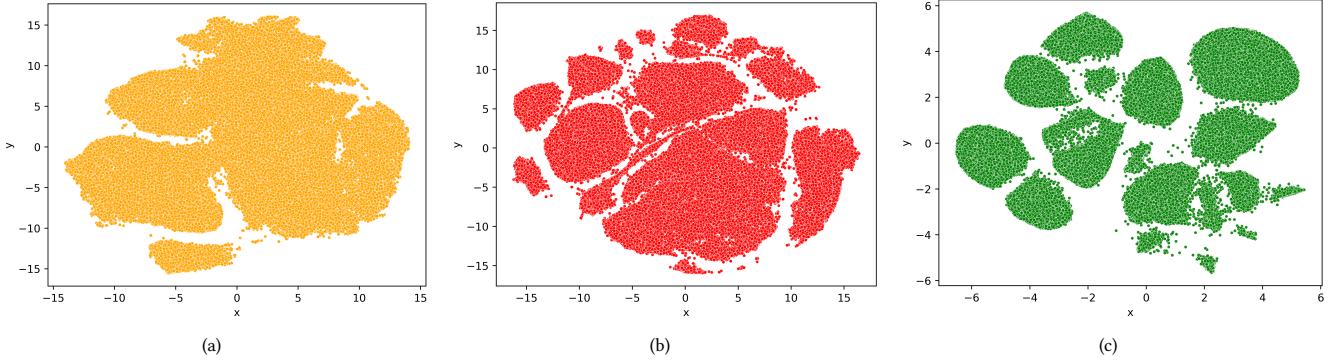
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXX.XXXXXXX>

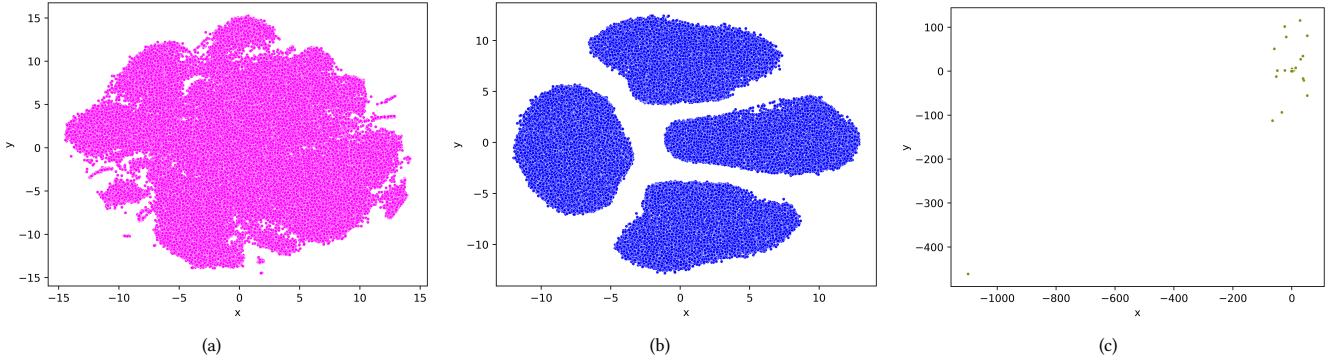
Deep features extracted from deep neural networks have been proven to capture the object representation well, no matter how small objects are or how diverse the dataset is [14], and recent object detectors built on CenterNet2[35], YOLOv4 [6], SOD [5], and TPH-YOLOv5 [37] provide an efficient 1D representation of 2D objects in a video. Similarity search looks for object representations in a database that are similar or close to a query based on a specific measure of similarity. That measure is usually a distance function. Let's define  $X$  as a metric space with associated distance function  $d(p, q)$ , and  $P$  as a set of points in that metric space  $p, q \in P$ . The *nearest neighbor* of a query point  $q$  is  $p$  if  $d(p, q) \leq d(q, p')$ ,  $\forall q, p, p' \in P, p \neq p'$ . The  $k$ -nearest neighbors ( $k$ -NN) search identifies the top  $k$  nearest neighbors to query  $q$  and has complexity  $O(|P| \times d)$ , where  $|P|$  is a number of points in  $P$  and  $d$  is the dimension of points in  $P$ . This approach does not scale for  $|P|$  in millions or billions and  $d$  in hundreds and thousands as the  $k$ -NN search considers the entire dataset each time a query is initiated. Our original task of finding a needle in the haystack and searching the entire haystack every time is reformulated as an indexing and search problem. First, we construct the data structure that summarizes point dataset  $P$  with the objective of enabling efficient nearest neighbor retrieval without the need to compute all distances from query vector  $q$  to all the points in  $P$  while trying to match the full exact retrieval set as closely as possible. The approximate nearest-neighbor (ANN) methods are traditionally optimized to balance efficiency and effectiveness and offer speed up at the account of the accuracy [1]. ANN methods are optimized to return *any* point  $p' \in P$  such that the distance from  $q$  to  $p'$  is at most  $c \cdot \min_{p \in P} D(q, p)$ , for some  $c \geq 1$  fast, and they can be roughly grouped as graph-based [11, 30, 33], hashing-based [20, 23, 34], and partition-based [3, 13] methods. The proposed methods can be applied to a variety of feature databases, and the effectiveness varies based on the size and application as indexing and search approaches are correlated with data characteristics. Deep descriptor databases are high dimensional ( 1024) and sparse ( 58% of the feature vectors are 0).

### 1.1 Motivation

The deep descriptor databases are high dimensional ( 1024+) and sparse as on average 58% of feature vector is 0. We summarized the comparisons with non-deep descriptor datasets in Table 1. Crawl840B is the word vector representation that has no 0 entries, and DEEP10M has only 19 entries with 0, while deep descriptor float databases DOTA2.0, DIOR, and Visdrone have on average 77% per feature vector and SIFT10M has 25% 0 entries as outlined in table 1. This motivated us to look into descriptor distribution along the dimension and their t-Sne visualizations. The data distributions along the first two dimensions for whole deep



**Figure 1: t-SNE distribution of (a)DOTA2.0, (b)DIOR, and (c)VisDrone (100,000 instances) show similar 2D space projection distribution.**



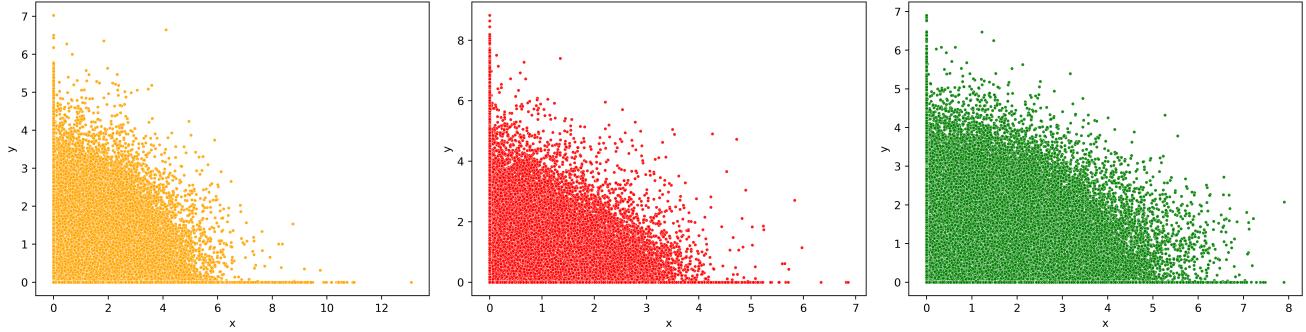
**Figure 2: t-SNE distribution for 100,000 features randomly sampled from (a) DEEP10M, (b) SIFT10M, and (c) Crawl840B datasets. The non-deep descriptor distribution differs in terms of cluster separation.**

float descriptor databases DOTA2.0, DIOR, and VisDrone in Figure 3 show that the data has the same distribution as the same DNN was used to extract 1024 dimensional vectors. The data has a vastly different distribution nature than the descriptor distribution of DEEP10M, SIFT10M, and Crawl840B in Figure 4. Note that the DEEP10M integer feature distribution differs from other deep descriptor databases because the actual output of 1024 dimensional deep feature of Googlenet's last fully connected layer [2] was compressed and normalized into 96-dimensional vectors using principal component analysis, and the Figure 4 illustrates the distribution of the first two principal components of the dataset. It is similar to the word vector of Crawl840B. The data distribution along the first two dimensions of SIFT10M image descriptor database shows that the distributions are skewed, with most of the values concentrated around a single value along each dimension. Therefore, the features are not uniformly distributed, and certain patterns are dominant in the feature space. This is because the SIFT algorithm [25] creates 16 patches from an image and utilizes the gradient orientation along 8 direction. Next, the t-SNE distributions of the first 100,000 vectors of sample three float databases (DOTA2.0, DIOR, and Visdrone) are shown in Figure 1, and the t-SNE distribution of sample integer descriptor database (DEEP10M, SIFT10M, and Crawl840B) are shown in Figure 2. While the Deep10M t-SNE visualization of the first 96 principal component vectors Figure 2(a) exhibits similarity to the overall mapping of the deep descriptors in Figure 1, SIFT descriptors

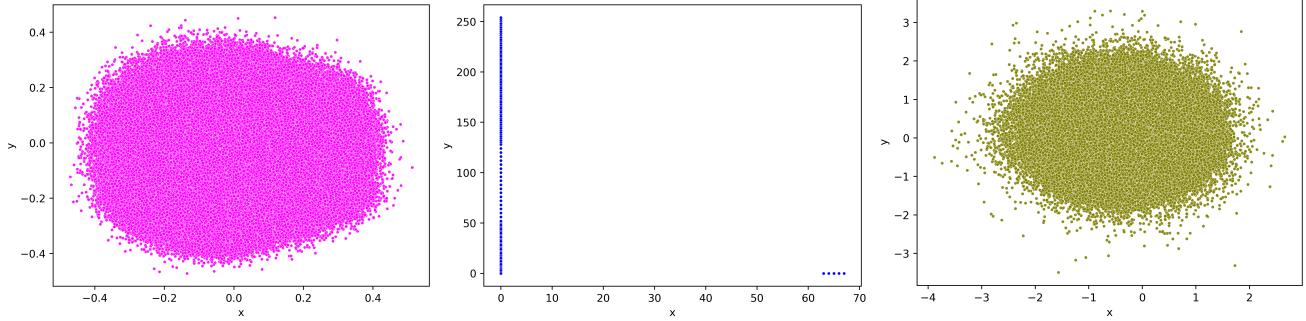
cluster well into 4 clusters. In contrast, word visualization fails as they all overlap in few dots in the top right corner in Figure 2(c). We identify a clear difference in behavior of deep descriptor databases and the compressed deep descriptor, traditional image descriptor, and word vector databases. Since our task at hand is to have the most effective and efficient way to find similar objects as represented by high dimensional sparse deep descriptors, we focus on proposing the indexing and search structure most suitable for that task.

## 1.2 Proposed Work

Our task focuses on effective and scalable indexing and retrieval improvements for real deep feature databases for unknown class discovery applications. We introduce the *Stratified Graph (SG)* as a solution for effective indexing and retrieval. The *Stratified Graph (SG)* approach is tailored to address the challenges of searching through high-dimensional and sparse datasets of deep descriptors by considering the sparsity of feature vectors. The Stratified Graph Indexing technique involves segmenting the data into layers, with each layer based on the distance of the point from the center of mass of the dataset. In each layer, the index is constructed as a bi-directional graph that links each vector to its closest neighbors within the same layer, and to its nearest neighbor in the adjacent outer layers. This layer graph is designed to achieve skip list properties which helps the searching algorithm to skip visiting nodes along the greedy search path [29]. During the search process, the



**Figure 3: Data distribution along first two dimensions for DOTA2.0, DIOR, and VisDrone.**



**Figure 4: Data distribution along first two dimensions for DEEP10M, SIFT10M and Crawl840B. Note that these are two first principal components of the Deep10M**

Stratified Graph search starts from the innermost layer and gradually moves outwards to the outer layers through the connected graph. The rest of the paper is organized as follows. In Section 2 we related work and describe the limitations of existing methods in terms of scalability and efficiency for larger datasets. In Section 3, we discuss the indexing and search procedure in detail. The advantage of the Stratified Graph is that it does not have a hierarchical structure and that SG connectivity at the same layer enhances the recall, or the ability to retrieve relevant results, while the SG connections between layers help maintain logarithmic complexity scaling for faster search speed. In Section 4, we present an in-depth comparison of multiple state-of-the-art methods in six data collections and compare the performance of the indexing and search methods in the word embedding, visual descriptor, and three deep descriptor databases in terms of high recall, precision, and F1-score at any depth of retrieval and fast retrieval times, and index size. We present our findings that the Stratified Graph approach is the most suitable for deep descriptor application in Section 5.

## 2 RELATED WORK

The *Fais* library [2] enables efficient partitioning of data in Voronoi cells [2], where the index of each cell is a centroid of that cell and product quantization is used [18] to compress data. The approach and its hierarchical improvements have not been shown to scale well for further retrieval results in 128 dimensional SIFT10M dataset with 10 million instances and 96 dimensional DEEP10M dataset with 10 million instances shown in Section 4. Annoy generates a number of hierarchical 2-means trees by recursive partitioning. Each iteration results in the formation of two centers by conducting a basic clustering algorithm on a subset of samples from the input

data points. The two centers define a partition hyperplane that is equidistant from each of them. The hyperplane then partitions the data points into two sub-trees, and the algorithm iteratively generates the index on each sub-tree [3]. This approach did not scale to larger high dimensional datasets in terms of the speed of recovery and accuracy of the retrieved results [24]. *Hierarchical Navigable Small World* (HNSW) [27] arranges the graph into a hierarchy of proximity graph layers with lower layer containing all the feature vectors and higher layers containing a subset of previous layers in the hierarchy. However, this architecture results in larger index size shown in Section 4. *Navigating Spread-out Graph* (NSG) [12] favors the “Navigating Node” to make the search efficient, but the sparsity of the data space and the indexing complexity does not scale well when the dimension of the feature vector grows [11]. *Navigating the satellite system graph* (NSSG) improves over SSG as it introduces the satellite system graph (SSG) and a more efficient pruning technique during index building to address the high-dimensional curse. The sparsity of NSSG can be controlled by a parameter, but the chances that the monotonic search stage fails are greater as the size and dimension of the database increases [11]. The *Tree-Based Search Graph* (TBSG) proposes to handle this problem with the probability of monotonic search success by combining the Cover Tree [4] and BKNNG (Bi-directed K-Nearest Neighbor Graph) [10] algorithms [9]. In the *Hierarchical Satellite System Graph* (HSSG), the nodes in the dataset are separated into layers, and NSGs are created on each layer separately. When searching in the high layer, the search process can skip a long distance, reducing the total number of steps in large data [33]. The index processes in separate layers are independent once the nodes are picked, and the index algorithm can be run distributively, decreasing the index algorithm’s time consumption.

During the search, HSSG performs a faster coarse search on the upper layer with fewer nodes. Following the coarse finds, HSSG conducts a more precise recursive fine search in the bottom layer at the cost of the high indexing and memory overhead compared to SSG [33]. *Neighborhood Graph Tree* (NGT) [15] uses a range search during the graph construction mechanism, and, to avoid a high degree of neighboring nodes and reduce memory overhead, applies a three-degree adjustment by connecting each feature vector to its three nearest neighbors throughout the graph. During the query process, NGT generates a seed using the VP tree [17] and performs a range search to obtain the nearest neighbors. A major drawback of NGT is that if the query and seed are far away from each other in the search space, then it takes many hops in between to reach the query from the seed, and thus increases the retrieval time. One way to address this problem is to transform the  $k$  nearest neighbor graph into a undirected one, and the other is to construct an undirected graph by continuously inserting elements [16]. Learned Index for large-scale **D**Ense passage **R**etrieval (LIDER)'s [31] hierarchical architecture is based on clustering and consists of two levels of core models. A core model is the basic unit of LIDER for indexing and searching data. It consists of an adapted Recursive Model Index (RMI) and a dimension reduction component that contains an expanded SortingKeys-LSH (SK-LSH) and a key re-scaling module. High-dimensional dense embeddings are converted into one-dimensional keys and sorted in a certain order to make quick predictions by the RMI. However, for a small number of clusters, each cluster yields a large number of feature vectors, making it more difficult for RMI to learn the distribution effectively. As a result, the quality of in-cluster retrieval degrades. On the contrary, for a large number of clusters, recall suffers. All neighboring ANN methods suffer from a long index-building time and low retrieval for large deep-descriptor databases [24]. In this paper we compare and contrast the proposed work with state-of-the-art for the large sparse descriptor retrieval.

### 3 STRATIFIED GRAPH METHODOLOGY

#### 3.1 Stratified Graph Indexing (SGI)

The Stratified Graph Indexing (SGI) arranges the feature vectors into layers based on their distances from the centroid. The center of the target is computed as the average of the sample or the entire dataset  $P$ . The feature vectors that are closer to the centroid are stratified in closer layers while the feature vectors that are farther away are stratified in farther layers. Figure 5(a) illustrates the Stratified Graph Indexing (SGI) for the Euclidean distance and the layers are stacked into the target shape in 2D. Note that for the Manhattan distance, the layers can be illustrated as a set of rotated squares in 2D illustration. The bidirectional graph is constructed by connecting each feature vector to its  $m$ -Nearest Neighbors. In Figure 5(a)  $m$  is set to 4. We specify the number of layers as  $\log_2 m + 1$ . Therefore, our example in Figure 5 has 3 layers. We select the layer width by dividing the distance to the farthest point from the center by the layer numbers. Therefore, all the layers have the same width. Four points in Layer 0 have two connections within the central layer and one each with the outside layers. In Layer 1, ten points have three connections within Layer 1 and one with the outer layer. All the points in Layer 2 have four connections within Layer 2. The subsequent layer connections allow the SGI algorithm to achieve

skip list properties which helps the Stratified Graph Retrieval (SGR) algorithm a faster search in the graph, which we discuss later in this section. Higher vertex degree  $m$ , higher index size, slower search but higher recall as each point is connected to its true  $m$  nearest neighbors (Figure 8).

The index-building phase involves determining the layers of each element based on their distances from the centroid and then constructing a kNN graph within each layer, from the outermost layer to the innermost layer. In this process, an outlier filtering factor  $f$  is used to filter out elements that are too far from the mean distance and ensure that outliers do not affect the layer boundaries.

---

#### Algorithm 1: Stratified Graph Indexing

---

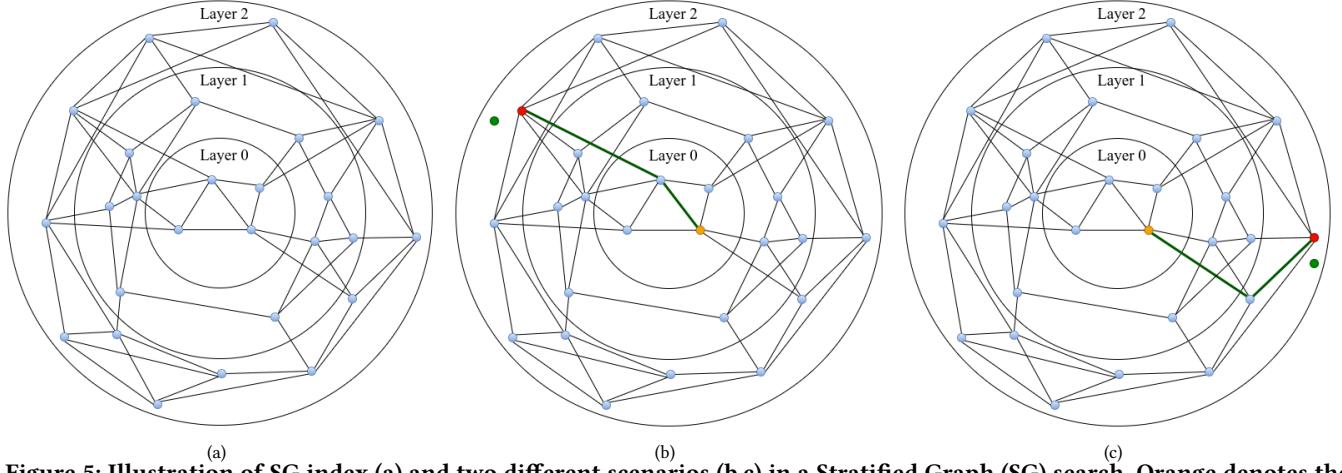
```

Input: stratified graph SG, dataset P, graph degree m,
        outlier filtering factor f
Output: (SG, X, M, cand, f)
1 graph  $\leftarrow \phi$ 
2 layerGraphList  $\leftarrow \phi$ 
3 layeredElem  $\leftarrow \text{LAYERING}(X, M, f)$ 
4 foreach maxLayer to minLayer of layeredElem do
5   foreach elem of layer do
6     clg  $\leftarrow \text{ADD}(\text{elem}, M, \text{cand})$ 
7     if layerGraphList not empty then
8       foreach g in layerGraphList do
9         n  $\leftarrow \text{SEARCH}(g, \text{elem}, k = 1, \text{cand} = 1)$ 
10        clg  $\leftarrow \text{update clg inserting } n \text{ to neighbor list}$ 
11        of elem
12      end
13    end
14    add clg to layerGraphList
15  end
16  m  $\leftarrow m - 1$ 
17 end
18 graph  $\leftarrow \text{merge all the graphs in layerGraphList}$ 

```

---

Algorithm 2 describes the process of determining the layers for each element. The algorithm computes the mean distance from the centroid,  $dist$ , and the standard deviation of distances,  $\sigma$ . Elements whose distances are greater than  $f$  times  $\sigma$  from the mean distance are filtered out. The remaining elements are then divided into layers based on their distances from the centroid, with each layer containing elements within a certain range of distances. After the layers are determined, the kNN graph is constructed for each layer (Algorithm 1 line 6). Then the algorithm adds the next-layer neighbors to all the feature vectors that are not in the outermost layer, which helps to capture the global structure of the data (Algorithm 1 line 7-12). Any new insertion of neighbors in the graph is simply an SG search (Algorithm 3) in the existing index. The final graph is constructed by taking a simple union of all the graphs for each layer (Algorithm 1 line 17). The resulting graph captures the local and global structure of the data and can be used for efficient similarity search.



**Figure 5: Illustration of SG index (a) and two different scenarios (b,c) in a Stratified Graph (SG) search. Orange denotes the starting point of a search, red denotes the nearest neighbor, and green edges show the path of the greedy algorithm to the query(shown in green).**

---

**Algorithm 2:** Graph Layering Method

---

**Input:** data vector  $X$ , number of established connections  $M$ , outlier filtering factor  $f$   
**Output:** LAYERING ( $X, M, f$ )  
1  $numLayer \leftarrow \log_2 M$   
2  $cen \leftarrow$  mean of  $X$   
3  $dist \leftarrow$  distances from centroid to all data vectors  
4  $avg \leftarrow$  mean of all distances  
5  $\sigma \leftarrow$  standard deviation of all distances  
6  $u_b \leftarrow avg + f \times \sigma$   
7  $l_b \leftarrow$  smallest of  $dist$   
8  $r \leftarrow \frac{u_b - l_b}{numLayer}$   
9  $layeredElem \leftarrow \emptyset$   
10 **foreach**  $(d, x)$  of  $(dist, X)$  **do**  
11   |  $l \leftarrow \frac{d}{r}$   
12   | add element  $x$  to layer  $l$  in  $layeredElem$   
13 **end**  
14 **return**  $layeredElem$

---

The stratified graph indexing (SG) has two phases. During the first phase, each element is added one at a time by iterative insertions, which are simply a series of ANN searches at different levels. Thus, the first phase has a complexity of  $O(|P|\log(|P|))$ .

The second phase of stratified graph (SG) index building is also a series of ANN searches at different layers. Thus, similar to the first phase, the second phase also has a complexity of  $O(|P|\log(|P|))$ . Therefore, the overall complexity of the index building of the stratified graph (SG) scales as  $O(|P|\log(|P|))$ .

### 3.2 The Stratified Graph Retrieval (SGR)

The Stratified Graph Retrieval (SGR) steps are illustrated in Figure 5(b) and (c). The search within an index starts with a random feature vector at the innermost layer denoted in orange in Figure 5(c) and

---

**Algorithm 3:** Stratified Graph Search

---

**Input:** graph index  $g$ , query element  $q$ , number of nearest neighbors  $k$ , size of dynamic candidate list  $cand$   
**Output:**  $k$  closest neighbors to  $q$   
1  $ep \leftarrow$  get entry point of  $g$   
2  $p \leftarrow$  extract nearest neighbor to  $q$  starting with  $ep$   
3  $C \leftarrow$  extract  $cand$  neighbors to  $p$  from  $g$   
4  $neighbors \leftarrow$  top  $k$  closest from  $C$  to  $q$   
5 **return**  $k$ -NN neighbors to  $q$

---

(d), where each feature has the highest number of next-layer neighbors. Then a greedy search within the graph is applied to retrieve the closest neighbor (denoted in red in Figure 5(c) and (d)) and top  $k$  to the query (denoted in green in Figure 5(c) and (d)) are returned: for this example,  $k$  is 1. A heap of size  $cand$  maintains the neighbor list based on their distances to the query along the search path. The parameter  $cand$  also determines the depth of search that can be performed in the graph. Layering enhances the search speed of the SGR algorithm by skipping visiting nodes if the current node and query are some layers apart from each other. Figure 5 (b) and Figure 5 (c) utilizing the next layer neighbor that is closer to the query in the feature space. In other cases, the search algorithm will perform a simple greedy search in the graph to retrieve the nearest neighbors to the query.

Algorithm 3 demonstrates how the greedy search process works in the Stratified Graph Retrieval (SGR) technique. Initially, the algorithm locates the nearest neighboring point  $p$  to the incoming query  $q$  by examining the neighbor list of the starting point  $ep$ . It then identifies the top  $k$  neighbors from  $p$  to query  $q$ . The number of hops and the average degree of the items on the greedy path is multiplied to obtain the total number of distance calculations. The SGR approach benefits from the outer layer connections that enable it to bypass visiting a considerable portion of the graph, leading to logarithmic time complexity. Therefore, the time complexity of the SGR technique can be expressed as  $O(\log(|P|))$ .

## 4 EXPERIMENTS

We measure the performance of our Stratified Graph (SG) method as compared to the six different state-of-the-art methods: Lightweight approximate Nearest-Neighbor library (N2) [21], Non-Metric Space Library (NMSlib) [7], Hierarchical Navigable Small World library (HNSWlib) [27], Facebook AI Similarity Search library (FaissHNSW) [19], Approximate Nearest-Neighbors Oh Yeah (Annoy) [3], and Hybrid Approximate Nearest-Neighbor Indexing and Search (HANNIS) library built on the HNSW algorithm [26].

**Table 1: Dataset characteristics**

Data	Desc. Type	Frame work	Dim ension	DB size in GB	# features	% os
VisDrone [36]	Video	[5]	1024	6.2	1.51	69
DOTA2.0 [32]	Image	[5]	1024	11.1	2.69	80
DIOR [22]	Image	[5]	1024	5.2	1.27	83
DEEP10M [2]	Image	[2]	96	3.8	10	0
SIFT10M [8]	Image	[8]	128	0.516	10	25
Crawl840B [28]	Text	[28]	300	5.6	2.2	0

### 4.1 Experimental Setup

Datasets characteristics used from the experiments are summarized in Table 1. **VisDrone** dataset contains 1,515,007 instances of 1024 dimensional object deep feature descriptors extracted from VisDrone video [36]. **DOTA2.0** dataset contains 2,697,873 instances of 1024-dimensional object deep feature descriptors extracted from DOTA2.0 [32]. **DIOR** dataset contains 1,278,863 instances of 1024 dimensional object deep feature descriptors extracted from DIOR [22]. For the VisDrone, DOTA2.0, and datasets, we extracted 1024 dimensional object-level deep features from the last fully-connected layer using pipeline SOD [5]. **DEEP10M** dataset contains 10 million instances of 96-dimensional floating vectors. The original 10 million 1024-dimensional image embedding outputs of the Googlenet's last fully connected layer [2] were compressed and normalized into 96-dimensional vectors using principal component analysis. **SIFT10M** dataset contains 10,000,000 instances of 128-dimensional integer SIFT image descriptors [25] extracted from Caltech-256 41  $\sqcap$  41 whole image patches [8]. Crawl840B dataset with 300 dimensions and 2.2 million instances of vector embeddings of common crawl words using GloVe [28]. Visdrone, DOTA2.0, DIOR, DEEP10M, SIFT10M, and Crawl840B dataset sizes are 6.2, 11.1, 5.2, 3.8, 5.6, 0.516, and 11.1, in Giga Bytes, respectively.

**Evaluation** of the indexing and search methods have been conducted from the task's perspective. Our task is focused on a visual search for unknown objects or class discovery in the petabytes of image and video archives. Therefore, the recall@ $k$  and precision@ $k$  need to remain consistent as  $k$  increases while keeping retrieval time and index size comparable to state of the art. Note that we include the word2Vec dataset and the SIFT10M databases to illustrate the varying effect of different methods based on the application. Four performance measurement metrics have been used to evaluate the performance of each method: recall@ $k$ , precision@ $k$ , retrieval time, and index size. The number  $k$  is the size of the retrieved set,  $k \in [5, 10, 20, 50, 100]$ . Let  $GT_k$  be the  $k$  Nearest Neighbors set of descriptor indices recovered by the brute force search, and let  $M_k$  be the set of size  $k$  recovered by the method  $M$ .  $|M_k \cap GT_k|$  is the

number of true positives, the number of nearest neighbors recovered by the method  $M$  in  $M_k$  that match  $GT_k$ . **Precision @ k** is measured as the fraction of the retrieval set that is relevant to the query:  $P_k = (|M_k \cap GT_k|)/|M_k|$ , and **Recall @ k** is measured as the fraction of the retrieval set that is relevant to the ground truth:  $R_k = (|M_k \cap GT_k|)/|GT_k|$ . We evaluate 7 methods  $M$ ,  $M \in [N2, NMSlib, HNSWlib, FaissHNSW, Annoy, HANNIS, SG]$ . **Retrieval time** defines the time to retrieve the nearest neighbor  $k$  for a single query in milliseconds. **Index size** defines the memory cost to save the indexes in the memory in Giga Bytes.

**Setup** All experiments were carried out on Ubuntu 20.04.3 server with 11th generation Intel® CoreTM i9-11900K @ 3.5GHzX16 CPU with 128GB RAM and NVIDIA GeForce RTX 3070 8GB mem GPU. The python implementation of SG library can be found in <https://anonymous.4open.science/r/SG-4644>.

### 4.2 Effectiveness of the Retrieval Methods

**Table 2: Precision and recall comparison for VisDrone.**

Metrics Methods	Precision					Recall				
	k=5	k=10	k=20	k=50	k=100	k=5	k=10	k=20	k=50	k=100
N2	0.24	0.22	0.14	0.09	0.04	0.56	0.56	0.61	0.66	0.69
NMSlib	0.34	0.18	0.09	0.04	0.02	0.48	0.41	0.38	0.32	0.32
HNSWlib	0.76	0.54	0.50	0.49	<b>0.38</b>	0.86	0.79	0.80	0.88	<b>0.91</b>
FaissHNSW	0.58	0.44	0.42	0.18	0.09	0.72	0.68	0.66	0.60	0.48
Annoy	0.40	0.27	0.21	0.19	0.13	0.50	0.50	0.53	0.66	0.73
HANNIS	0.66	0.55	0.43	0.32	0.26	0.74	0.75	0.74	0.80	0.80
SG	<b>0.80</b>	<b>0.79</b>	<b>0.74</b>	<b>0.61</b>	0.36	<b>0.96</b>	<b>0.95</b>	<b>0.96</b>	<b>0.94</b>	<b>0.91</b>

**Table 3: Precision and recall comparison for DOTA2.0.**

Metrics Methods	Precision					Recall				
	k=5	k=10	k=20	k=50	k=100	k=5	k=10	k=20	k=50	k=100
N2	0.80	0.88	0.83	<b>0.80</b>	0.66	0.86	0.97	0.97	<b>0.97</b>	<b>0.97</b>
NMSlib	0.28	0.22	0.13	0.05	0.03	0.58	0.71	0.67	0.62	0.60
HNSWlib	0.96	0.79	0.66	0.65	0.44	0.98	0.93	0.95	0.95	0.95
FaissHNSW	0.71	0.60	0.46	0.20	0.09	0.84	0.83	0.78	0.74	0.62
Annoy	0.69	0.48	0.32	0.39	0.33	0.86	0.78	0.75	0.82	0.85
HANNIS	0.94	0.88	0.83	0.75	<b>0.70</b>	0.98	0.98	0.98	0.96	0.96
SG	<b>1</b>	<b>1</b>	<b>0.98</b>	0.72	0.51	<b>1</b>	<b>1</b>	<b>0.99</b>	<b>0.97</b>	0.95

**Table 4: Precision and recall comparison for DIOR.**

Metrics Methods	Precision					Recall				
	k=5	k=10	k=20	k=50	k=100	k=5	k=10	k=20	k=50	k=100
N2	<b>1</b>	0.85	0.78	0.73	0.71	<b>1</b>	0.97	0.98	0.99	<b>0.99</b>
NMSlib	0.27	0.18	0.11	0.05	0.02	0.66	0.74	0.76	0.73	0.73
HNSWlib	0.93	0.91	0.95	0.89	0.90	0.96	0.98	0.99	0.99	0.99
FaissHNSW	0.9	0.9	0.79	0.37	0.16	0.9	0.9	0.87	0.86	0.72
Annoy	0.71	0.63	0.62	0.52	0.54	0.88	0.87	0.89	0.92	0.94
HANNIS	<b>1</b>	<b>1</b>	<b>1</b>	0.98	0.90	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>0.99</b>
SG	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>0.93</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>0.99</b>

In this experiment, we compare the Stratified Graph approach with six existing approaches for 2.7 million DOTA 2.0, 1.3 DIOR and 1.5 million VisDrone 1024 dimensional databases created using the methodology published in [5]. Table 3, 4 and 2 shows that Stratified Graph (SG) is the most suitable algorithm to match deep features and thus uncover similar unlabeled objects for the DOTA2.0, DIOR, and VisDrone datasets in terms of precision and recall. N2 and HNSWlib perform on par with Stratified Graph (SG) and in terms of effectiveness *only* for larger retrieval sets as illustrated

in Table 3, 4 and 2. FaissHNSW and Annoy performance quickly degrades for  $k > 5$  so the methods are not suitable for deep descriptor database matching. NMSlib has consistently low precision and low recall for all  $k$  in Table 3, 4 and 2.

**Table 5: Precision and recall comparison for DEEP10M.**

Metrics Methods	Precision					Recall				
	k=5	k=10	k=20	k=50	k=100	k=5	k=10	k=20	k=50	k=100
N2	0.30	0.30	0.30	0.21	0.27	0.62	0.75	0.79	<b>0.84</b>	<b>0.87</b>
NMSlib	<b>0.96</b>	<b>0.81</b>	0.60	0.25	0.12	<b>0.98</b>	<b>0.95</b>	<b>0.90</b>	0.76	0.65
HNSWlib	0.70	0.56	0.46	0.33	0.30	0.80	0.75	0.77	0.80	0.85
FaissHNSW	0.69	0.45	0.29	0.11	0.05	0.78	0.69	0.60	0.56	0.45
Annoy	0.30	0.23	0.13	0.22	0.19	0.54	0.58	0.64	0.73	0.78
HANNIS	0.76	0.61	0.49	0.38	0.34	0.86	0.80	0.79	0.80	<b>0.87</b>
SG	0.78	0.72	<b>0.63</b>	<b>0.52</b>	<b>0.41</b>	0.92	0.88	0.87	<b>0.84</b>	0.84

**Table 6: Precision and recall comparison for SIFT10M.**

Metrics Methods	Precision					Recall				
	k=5	k=10	k=20	k=50	k=100	k=5	k=10	k=20	k=50	k=100
N2	0.56	0.47	0.36	0.24	0.18	0.78	0.82	0.82	0.86	0.87
NMSlib	0.22	0.12	0.06	0.03	0.01	0.64	0.66	0.63	0.58	0.52
HNSWlib	0.85	0.76	0.52	0.27	0.25	0.94	0.91	0.84	0.82	0.84
FaissHNSW	0.71	0.43	0.24	0.08	0.04	0.86	0.79	0.72	0.58	0.46
Annoy	0.13	0.21	0.16	0.09	0.11	0.42	0.44	0.50	0.60	0.71
HANNIS	0.86	0.74	0.63	<b>0.42</b>	0.22	0.94	0.93	0.93	<b>0.92</b>	<b>0.90</b>
SG	<b>0.94</b>	<b>0.83</b>	<b>0.74</b>	<b>0.42</b>	<b>0.27</b>	<b>0.96</b>	<b>0.98</b>	<b>0.96</b>	0.90	0.82

**Table 7: Precision and recall comparison for Crawl840B.**

Metrics Methods	Precision					Recall				
	k=5	k=10	k=20	k=50	k=100	k=5	k=10	k=20	k=50	k=100
N2	0.61	0.65	0.65	0.65	0.48	0.66	0.78	0.91	0.95	<b>0.97</b>
NMSlib	0.55	0.40	0.26	0.11	0.05	0.78	0.78	0.72	0.64	0.57
HNSWlib	0.2	0.15	0.16	0.18	0.25	0.2	0.25	0.36	0.53	0.67
FaissHNSW	0.3	0.25	0.16	0.06	0.03	0.42	0.45	0.45	0.38	0.28
Annoy	0.53	0.50	0.39	0.35	0.34	0.76	0.76	0.76	0.76	0.74
HANNIS	<b>0.94</b>	<b>0.92</b>	<b>0.91</b>	<b>0.83</b>	<b>0.76</b>	<b>0.98</b>	<b>0.97</b>	<b>0.98</b>	<b>0.96</b>	0.96
SG	70	0.64	0.57	0.41	0.31	0.70	0.70	0.70	0.69	0.70

**Visual retrieval** results for a single query for DIOR dataset are shown in Figure 6. We compare retrieved images for 10-NN search with respect to the brute force result shown as ground truth in Figure 6. All the methods return similar top 4 results as brute force results. Therefore, we only show the fifth and sixth retrieval results for all the methods in 6. For this particular object deep descriptor query, N2, HANNIS and SG retrieves **all 6** unique images correctly. Annoy retrieves 4 out of 6, and the rest retrieves 5 out of 6 images correctly. Some queries may be difficult for one method but easy for another, so we average all results throughout the experiment for 10 distinct queries.

Next, let's see if the methods behave in a similar fashion to different DNN used to extract descriptors. **DEEP10M** dataset contains 10 million instances of 96-dimensional floating vectors. The original 10 million 1024-dimensional image embedding outputs of the Googlenet's last fully connected layer [2] were compressed and normalized into 96-dimensional vectors using principal component analysis. Precision@ $k$  and recall@ $k$ -retrievals in Table 5 demonstrate that the Stratified Graph (SG) is performing well in the effectiveness of retrieval at higher  $k$ . The three best competitors of SG are N2, NMSlib, and HANNIS show inconsistent retrieval performance in Table 5. We also observe an interesting behavior of

N2 and Annoy in Table 5: the effectiveness of the indexing method is *improving* with larger  $k$ . HNSWlib performs moderately, and FaissHNSW has consistently poor performance in Table 5 compared to SG. Our interpretation from Table 5 result is: though DEEP10M features were extracted from DNN, the compression with principle component analysis alters the original characteristics of the dataset. However, Stratified Graph (SG) is shown to be the most consistent and suitable algorithm for discovering unknown classes for the DEEP10M dataset. In summary, SG is shown to be most robust algorithm for discovering unknown classes in deep descriptor databases.

Table 6 shows the precision@ $k$  and recall@ $k$  retrieval results for the SIFT10M dataset for SG and six comparing methods. Here, SG shows the dominating performance in precision@ $k$  retrieval results at **all**  $k \in [5, 10, 20, 50, 100]$  over the comparing methods. HANNIS is shown to be the best competitor of SG in recall@ $k$  for higher retrieval results in Table 6. N2 and Annoy show an upward trend in Recall@ $k$  retrieval with larger values of  $k$  in Table 6. FaissHNSW and NMSlib show consistently low performance than SG for both precision@ $k$  and recall@ $k$  retrieval results in Table 6. HNSWlib performs well and achieves similar precision@ $k$  and recall@ $k$  for higher retrieval results.

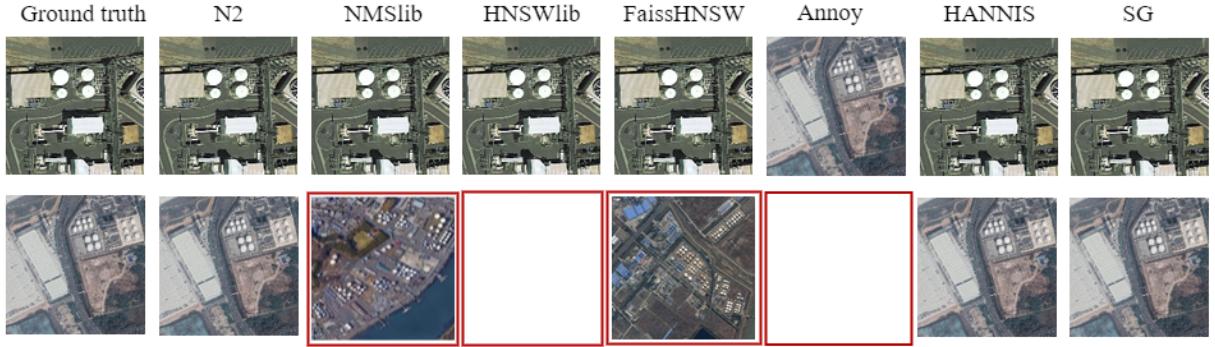
The precision@ $k$  and recall@ $k$  retrieval results for Crawl840B dataset are shown in Table 7. Though our algorithm is specifically designed for deep descriptor database, the performance of SG in Table 7 is compatible with the comparing methods for the vector representation of word embeddings data Crawl840B.

In summary, though SG is specifically designed for similarity search over deep descriptors, its performance is compatible with state-of-the-art algorithms for other descriptor databases.

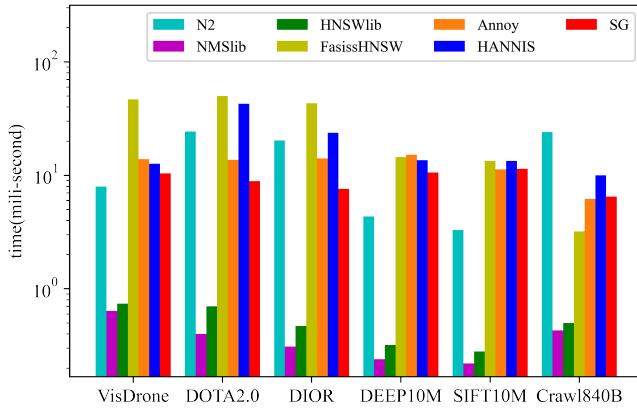
### 4.3 Efficiency Comparisons of the Index Size and Retrieval Times

In this experiment, we analyze the seven methods' retrieval time on a logarithmic scale for the six datasets *per method*. The Stratified Graph (SG) library is written in python without any optimization for speed. However, the retrieval times for SG in Figure 7 shows the promising result on 100 nearest neighbor search. Stratified Graph Retrieval (SGR) is faster than FaissHNSW, Annoy and HANNIS library for all six datasets except Crawl840B in Figure 7. Moreover, SGR is faster than N2 for DOTA2.0, DIOR and Crawl840B datasets. NMSlib and HNSWlib are faster than SGR for all six datasets. The retrieval time per method (7) for  $k = 100$  shows N2, NMSlib, HNSWlib, FaissHNSW, and HANNIS, the retrieval time corresponds to the dimension of the dataset except for the Crawl840B dataset. Annoy has moderate retrieval time and does not correlate with the dataset size in instances and the dimension and type of feature.

The Stratified Graph (SG) seems to do well with larger datasets. Table 8 shows the memory cost of saving the indexes in the memory. SG has the smallest index sizes than the comparing methods for all six datasets other than Crawl840B. NMSlib, HNSWlib, and FaissHNSW have similar index sizes in the memory. Annoy and HANNIS have larger index sizes than all the comparing methods. N2, NMSlib, HNSWlib, FaissHNSW, and HANNIS algorithms are built on HNSW algorithm. HNSW arranges the feature vectors in a hierarchical layer of proximity graphs where the upper layers in the



**Figure 6: Fifth and sixth retrieval results for a single query w.r.t. brute force search for seven methods N2, NMSlib, FaissHNSW, Annoy, HNSWlib, HANNIS, and Stratified Graph (SG).**



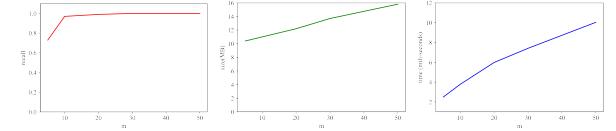
**Figure 7: Retrieval time for 100 nearest neighbor search per dataset for 7 approaches on 6 datasets.**

**Table 8: Index sizes (GB) per dataset for 7 approaches using 75 trees for Annoy and 16 neighbor connections for the rest of the six methods.**

Methods/ Datasets	N2	NMSlib	HNSWlib	FaissHNSW	Annoy	HANNIS	SG
DEEP10M	4.3	5.4	5.3	5.3	12.8	8.7	<b>2.3</b>
SIFT10M	2.2	2.7	2.6	2.6	5.6	3.9	<b>1.1</b>
Crawl840B	<b>2.7</b>	3	3	3	4.6	3.7	2.8
DOTA2.0	11.2	11.5	11.5	11.4	14.2	12.3	<b>3.2</b>
DIOR	5.3	5.3	5.4	5.4	6.8	5.8	<b>1.4</b>
VisDrone	6.2	6.2	6.4	6.4	7.6	6.9	<b>2.5</b>

hierarchy are subsets of the lower layer. Therefore, the graph index containing the proximity graphs has to store way more edge lists than SG, resulting in a higher memory overhead. Annoy has the largest index size of all the comparing methods because it requires storing many trees for better performance. Overall, SG requires up to four times less memory than comparing methods.

#### 4.4 Experiment 3: Ablation Study for the number of neighbors $m$ in SG



**Figure 8: Demonstration of recall, index size, and retrieval time increasing with a higher degree of  $m$ .**

The suitable value of neighbor connection  $m$  during index building depends on the characteristics of the deep descriptor databases and can range from 5 to 48, and the ablation study for random dataset of dimension  $x$  and size  $y$  is illustrated in Figure 8. The recall increases and approaches 1 as we increase the number of neighbor connections  $m$  from 5 to 48 in Figure 8. For this random dataset, we achieve a recall of 1 at around  $m = 16$ . We see a similar trend for index size and retrieval time, both increasing with the value of  $m$ . Therefore, the trade-off between effectiveness and efficiency depends on the number of neighbors connection  $m$ . Lower value of  $m$  results in faster search loosing some accuracy.

## 5 CONCLUSION

The potential hidden within video archives is immense, but the challenge lies in the limited amount of annotated images and objects. The solution is to identify objects similar to one another in feature space, which can be accomplished through a similarity search in deep descriptors databases. We propose Stratified Graph (SG) indexing and search as an effective solution for deep-descriptor matching in large, diverse databases defined for multiple real sets. Furthermore, the proposed stratified graph (SG) method outperforms state-of-the-art indexing and searching approaches in terms of recall and precision at the cost of slightly higher retrieval times. The precision, and recall improve up to 8% at depth 100 for the deep feature databases. Moreover, SG reduces the memory cost up to **four** time than comparing methods. As the next step, we plan to optimize the Stratified Graph (SG) method for efficient deep descriptor matching in billion deep descriptor databases.

## REFERENCES

- [1] Sunil Arya, David M Mount, Nathan S Netanyahu, Ruth Silverman, and Angela Y Wu. 1998. An optimal algorithm for approximate nearest neighbor searching fixed dimensions. *Journal of the ACM (JACM)* 45, 6 (1998), 891–923.
- [2] Dmitry Baranichuk, Artem Babenko, and Yury Malkov. 2018. Revisiting the Inverted Indices for Billion-Scale Approximate Nearest Neighbors. *CoRR* abs/1802.02422 (2018). arXiv:1802.02422 <http://arxiv.org/abs/1802.02422>
- [3] Erik Bernhardsson. 2018. *Annoy: Approximate Nearest Neighbors in C++/Python*. <https://pypi.org/project/annoy/> Python package version 1.17.1.
- [4] Alina Beygelzimer, Sham Kakade, and John Langford. 2006. Cover trees for nearest neighbor. In *Proceedings of the 23rd international conference on Machine learning*, 97–104.
- [5] Debojyoti Biswas and Jelena Tešić. 2023. Progressive Domain Adaptation with Contrastive Learning for Object Detection in the Satellite Imagery. arXiv:2209.02564 [cs.CV]
- [6] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. 2020. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934* (2020).
- [7] Leonid Boytsov and Bilegsaikhan Naidan. 2013. Engineering efficient and effective non-metric space library. In *International Conference on Similarity Search and Applications*. Springer, 280–293.
- [8] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>
- [9] Xiaobin Fan, Xiaoping Wang, Kai Lu, Lei Xue, and Jinjing Zhao. 2022. Tree-based Search Graph for Approximate Nearest Neighbor Search. *arXiv preprint arXiv:2201.03237* (2022).
- [10] Cong Fu and Deng Cai. 2016. Efanna: An extremely fast approximate nearest neighbor search algorithm based on knn graph. *arXiv preprint arXiv:1609.07228* (2016).
- [11] Cong Fu, Changxu Wang, and Deng Cai. 2021. High dimensional similarity search with satellite system graph: Efficiency, scalability, and unindexed query compatibility. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [12] Cong Fu, Chao Xiang, Changxu Wang, and Deng Cai. 2017. Fast approximate nearest neighbor search with the navigating spreading-out graph. *arXiv preprint arXiv:1707.00143* (2017).
- [13] Antonio Javier Gallego, Juan Ramón Rico-Juan, and Jose J Valero-Mas. 2022. Efficient k-nearest neighbor search based on clustering and adaptive k values. *Pattern recognition* 122 (2022), 108356.
- [14] David Heyse, Nicholas Warren, and Jelena Tešić. 2019. Identifying maritime vessels at multiple levels of descriptions using deep features. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, Tien Pham (Ed.), Vol. 11006. International Society for Optics and Photonics, SPIE, 423 – 431. <https://doi.org/10.1117/12.2519248>
- [15] Masajiro Iwasaki. 2015. Ngt: Neighborhood graph and tree for indexing.
- [16] Masajiro Iwasaki. 2016. Pruned bi-directed k-nearest neighbor graph for proximity search. In *International Conference on Similarity Search and Applications*. Springer, 20–33.
- [17] Masajiro Iwasaki and Daisuke Miyazaki. 2018. Optimization of indexing based on k-nearest neighbor graph for proximity search in high-dimensional data. *arXiv preprint arXiv:1810.07355* (2018).
- [18] Herve Jegou, Matthijs Douze, and Cordelia Schmid. 2010. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence* 33, 1 (2010), 117–128.
- [19] Jeff Johnson, Matthijs Douze, and Hervé Jegou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.
- [20] Sunwoo Kim, Haici Yang, and Minje Kim. 2020. Boosted locality sensitive hashing: Discriminative binary codes for source separation. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 106–110.
- [21] GeonHee Lee. 2017. TOROS N2 - lightweight approximate Nearest Neighbor library which runs fast even with large datasets. <https://github.com/kakao/n2> Python package version 0.1.7.
- [22] Ke Li, Gang Wan, Gong Cheng, Liqiu Meng, and Junwei Han. 2020. Object detection in optical remote sensing images: A survey and a new benchmark. *ISPRS Journal of Photogrammetry and Remote Sensing* 159 (2020), 296–307.
- [23] 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. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*. 289–300. <https://doi.org/10.1109/ICDE48307.2020.00032>
- [24] Wen Li, Ying Zhang, Yifang Sun, Wei Wang, Mingjie Li, Wenjie Zhang, and Xuemin Lin. 2019. Approximate nearest neighbor search on high dimensional data: experiments, analyses, and improvement. *IEEE Transactions on Knowledge and Data Engineering* 32, 8 (2019), 1475–1488.
- [25] G Lowe. 2004. Sift-the scale-invariant feature transform. *Int. J. 2*, 91-110 (2004), 2.
- [26] M M Mahabubur Rahman and Jelena Tešić. 2022. Hybrid Approximate Nearest Neighbor Indexing and Search (HANNIS) for Large Descriptor Databases. In *2022 IEEE International Conference on Big Data (Big Data)*. 3895–3902. <https://doi.org/10.1109/BigData55660.2022.10020464>
- [27] Ya Malkov and Dmitry A Yashunin. 2018. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence* 42, 4 (2018), 824–836.
- [28] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 1532–1543.
- [29] William Pugh. 1990. Skip lists: a probabilistic alternative to balanced trees. *Commun. ACM* 33, 6 (1990), 668–676.
- [30] M M Mahabubur Rahman and Jelena Tešić. 2022. Evaluating Hybrid Approximate Nearest Neighbor Indexing and Search (HANNIS) for High-dimensional Image Feature Search. In *2022 IEEE International Conference on Big Data (Big Data)*. 6802–6804. <https://doi.org/10.1109/BigData55660.2022.10021048>
- [31] Yifan Wang, Haodi Ma, and Daisy Zhe Wang. 2022. LIDER: An Efficient High-dimensional Learned Index for Large-scale Dense Passage Retrieval. *arXiv preprint arXiv:2205.00970* (2022).
- [32] Gui-Song Xia, Xiang Bai, Jian Ding, Zhen Zhu, Serge Belongie, Jiebo Luo, Mihai Datcu, Marcelli Pelillo, and Liangpei Zhang. 2018. DOTA: A large-scale dataset for object detection in aerial images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3974–3983.
- [33] Jiaru Zhang, Ruhui Ma, Tao Song, Yang Hua, Zhengui Xue, Chenyang Guan, and Haibing Guan. 2022. Hierarchical Satellite System Graph for Approximate Nearest Neighbor Search on Big Data. *ACM/IMS Transactions on Data Science (TDS)* 2, 4 (2022), 1–15.
- [34] Bolong Zheng, Zhao Xi, 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 (2020), 643–655.
- [35] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl. 2021. Probabilistic two-stage detection. *arXiv preprint arXiv:2103.07461* (2021).
- [36] Pengfei Zhu, Longyin Wen, Xiao Bian, Haibin Ling, and Qinghua Hu. 2018. Vision meets drones: A challenge. *arXiv preprint arXiv:1804.07437* (2018).
- [37] Xingkui Zhu, Shuchang Lyu, Xu Wang, and Qi Zhao. 2021. TPH-YOLOv5: Improved YOLOv5 based on transformer prediction head for object detection on drone-captured scenarios. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2778–2788.