

Stratified Graph Indexing for Efficient Search in Deep Descriptor Databases

M M Mahabubur Rahman¹ and Jelena Tešić^{1*}

¹Computer Science, Texas State University, 601 University Dr, San Marcos, 78666, Texas, USA.

*Corresponding author(s). E-mail(s): jtesic@txstate.edu;
Contributing authors: toufik@txstate.edu;

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. 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 enhance Average Recall (AR), 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 SG indexing and search has smaller memory usage (up to four times) and further Mean Average Precision and AR improves up to 8% than state-of-art for all six datasets at five retrieval depths.

Keywords: Similarity search, High-dimensional indexing, Deep descriptors search, Information retrieval, Vector search, Graph-based index

1 Introduction

Video archives house immense data and analytics that cannot be processed manually. The task of identifying similar objects in exabytes of video archive data in a scalable and effective way can help resolve many video analytics tasks, for example, the surveillance task: *When was this object spotted before?*; the crowd sensing task *Was this object in the archive footage?*, and the reconnaissance task: *Find me similar objects in archived footage*.

Classifiers sparsely help as object labeling is sparse, and labels are unavailable for every target of interest. Deep descriptors capture objectness characteristics well for the long tail of tasks, regardless of the size or density of the objects, as shown in [1, 2] for overhead imagery. Recent object detectors built on CenterNet2[3], YOLOv4 [4], SOD [5], and TPH-YOLOv5 [6] capture the objectness in a feature vector well. Now, identifying similar objects in exabytes of video archives reduces to designing an efficient indexing and search system in the deep descriptor databases [7]. This research proposes a novel indexing and search approach that supports efficient and effective *Approximate Nearest Neighbor* search in high-dimensional deep descriptor space.

1.1 Motivation

First, deep descriptor databases are high dimensional (1024) and sparse: for example, 58% of the feature vectors are 0 in the *TODO: which – Toufik* dataset. Thus, the sparsity and high dimensionality curse will degrade the similarity search performance compared to traditional descriptor databases with a low percentage of zero entries and smaller dimensions. Table 1 compares deep descriptor databases with non-deep descriptor databases: Crawl840B dataset contains no zero entries, DEEP10M has only 19 entries with zeros, and SIFT10M has only a quarter of values 0. Deep descriptors databases, extracted using [2] on DOTA2.0, DIOR, and Visdrone repositories, contain, on average, 77% zeros per feature vector.

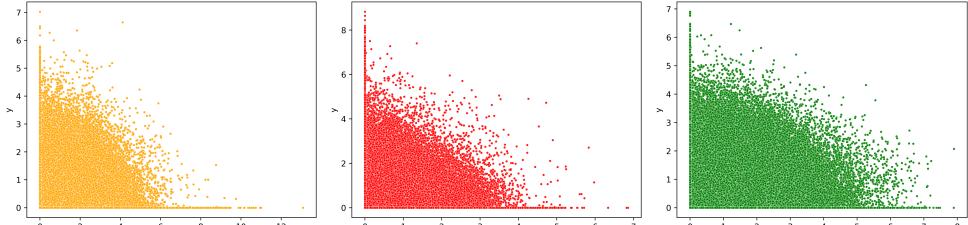


Fig. 1 Deep descriptor values for the first two^x dimensions for VisDrone, DOTA2.0, and DIOR datasets. *TODO: what is x, what is y? Missing labels. – Toufik*

Second, deep neural network architecture produces descriptor databases with similar distributions over domain datasets, as illustrated in Figure 1. The distribution differs from the two first principal components of Deep10M, SIFT feature extractor and Crawl840B data and feature extractor as shown in Figure 2. DEEP10 M's integer feature distribution differs from other deep descriptor databases because it compresses and normalizes the 1024-dimensional deep feature vectors from Googlenet's last fully connected layer into 96-dimensional vectors using principal component analysis. Figure 2 illustrates the distribution of this dataset's first two main components, which is somewhat similar to the word vector distribution of Crawl840B. The data distribution along the first two dimensions of the SIFT10M image descriptor database reveals a skewed distribution, with values concentrated around a single point along each dimension. This indicates non-uniform distribution, with specific dominant patterns in the feature space, which can be attributed to the SIFT [8] algorithm's approach of creating 16 patches from an image and using gradient orientation along eight directions. *TODO: Ideally, you will have several deep descriptors run on the same dataset*

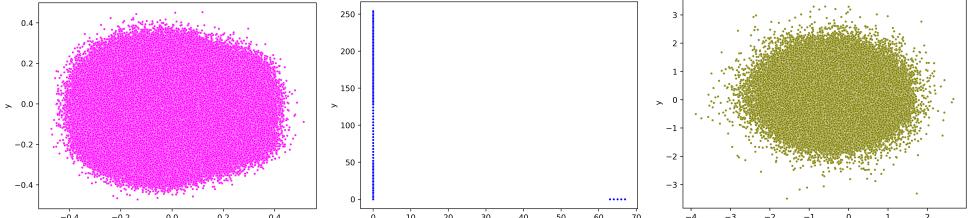


Fig. 2 Data distribution along first two dimensions for DEEP10M, SIFT10M and Crawl840B. Note that these are the two first principal components of the Deep10M

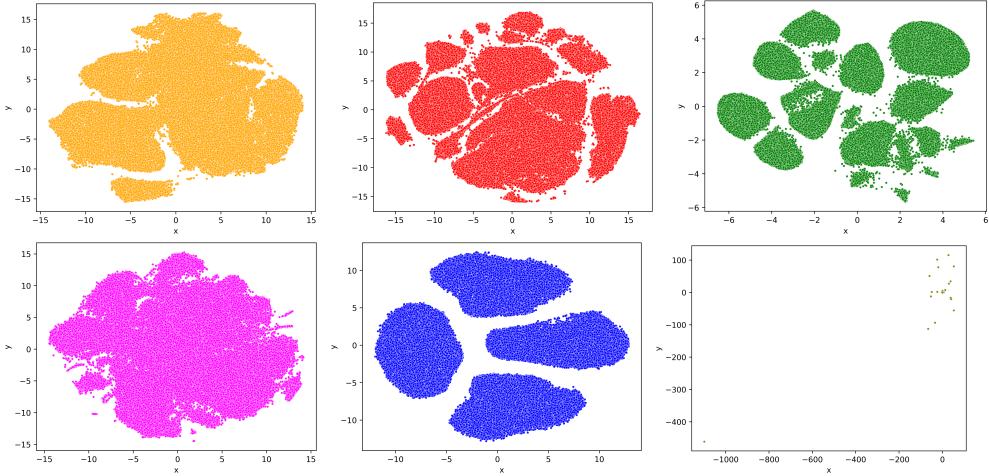


Fig. 3 t-SNE distribution for 100,000 features randomly sampled from (a) VisDrone, (b) DOTA2.0, (c) DIOR, (d) DEEP10M, (e) SIFT10M, and (f) Crawl840B datasets. *TODO: we are comparing apples and oranges here ... we need comparison of different DNN extractors and non-DNN extractors on the same dataset. – Toufik*

and compare them in Figures 1 and 2. This comparison makes no sense as in one figure, you are running the same system on three datasets; in the other, you have three methods on 3 different datasets. This has to change for the submission. – Toufik

Further insights come from the t-SNE visualizations. In Figure 3(a,b,c), we present t-SNE distributions of random 100,000 vectors from three float databases (Visdrone, DOTA2.0 and DIOR), while Figure 3(d,e,f) displays the t-SNE distribution of a sample of integer descriptor databases (DEEP10M, SIFT10M, and Crawl840B). Notably, the t-SNE visualization of DEEP10 M’s first 96 principal component vectors (Figure 3(d)) bears similarity to the overall mapping of the deep descriptors seen in Figure 3(a,b,c). However, SIFT descriptors cluster distinctly into four clusters (Figure 3(e)). In contrast, word vector visualization encounters challenges, with most points overlapping in a few dots in the top right corner (Figure 3(f)). Our investigation reveals a marked difference in the behavior of deep descriptor databases compared to compressed deep descriptors, traditional image descriptors, and word vector databases. Given our objective of effectively and efficiently finding similar objects represented by high-dimensional sparse deep descriptors, our focus is on proposing an indexing and search structure suited to this task.

2 State Of The Art

Approximate Nearest Neighbor Indexing and Search

Similarity search looks for object representations in a database that are similar or close to a query based on a specific similarity measure. 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 to enable 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 [9]. 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 $c \geq 1$. In Figure 4, the distance to the nearest neighbor $\min_{p \in P} D(q, p)$ from query q is shown as r . Therefore, the approximate nearest neighbor p' from query q is at most within distance $c \cdot r$.

Related Work

State of the art can be roughly grouped as graph-based [10–12], hashing-based [13–15], and partition-based [16, 17] methods. The *Faiss* library [18] enables efficient partitioning of data in Voronoi cells [18], where the index of each cell is a centroid of that cell and product quantization is used [19] to compress data. Annoy generates several hierarchical 2-means trees by recursive partitioning. Each iteration forms 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 other. The hyperplane then partitions the data points into two sub-trees, and the algorithm iteratively generates the index on each sub-tree [16]. This approach did not scale to larger high dimensional datasets in terms of the speed of recovery and accuracy of the retrieved results [20]. *Hierarchical Navigable Small World* (HNSW) [21] arranges the graph into a hierarchy of proximity graph layers with a lower layer containing all the feature vectors and higher layers containing a subset of previous layers in the hierarchy. However, this architecture results in a larger index size shown in Section 5. *Navigating Spread-out Graph* (NSG) [22] favors the “Navigating Node” to make the

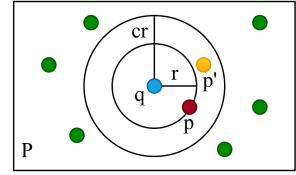


Fig. 4 Illustration of Approximate Nearest Neighbor search where p (red) denotes the nearest neighbor and p' (orange) denotes the approximate nearest neighbor with respect to query q (blue) for dataset P .

search efficient. Still, the sparsity of the data space and the indexing complexity do 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. A parameter can control the sparsity of NSSG, but the chances that the monotonic search stage fails are more significant 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 [23] and BKNNG (Bi-directed K-Nearest Neighbor Graph) [24] algorithms [25]. 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 number of steps in extensive data [12]. 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. HSSG performs a faster coarse search on the upper layer with fewer nodes during the search. Following the coarse finds, HSSG conducts a more precise recursive OK search in the bottom layer at the cost of the high indexing and memory overhead compared to SSG [12]. *Neighborhood Graph Tree* (NGT) [26] 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 [27] and performs a range search to obtain the nearest neighbors. A significant drawback of NGT is that if the query and seed are far from each other in the search space, it takes many hops between to reach the query from the root, and thus increases the retrieval time. One way to address this problem is to transform the k nearest neighbor graph into an undirected one, and the other is to construct an undirected graph by continuously inserting elements [28]. *Learned Index for large-scale DENSE passage Retrieval* (LIDER)'s [29] 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 critical re-scaling module. High-dimensional dense embeddings are converted into one-dimensional keys and sorted to make quick predictions by the RMI. However, for a small number of clusters, each cluster yields many 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 return for large deep-descriptor databases [20]. This paper compares and contrasts the proposed work with state-of-the-art for the large sparse descriptor retrieval.

TODO: Add ACM 2023 and special IEEE papers here. – Toufik

Stratified Graph Approach

We propose the *Stratified Graph* (SG) indexing (SGI) and retrieval (SGR) algorithms to address the challenges of searching through high-dimensional and sparse datasets of

deep descriptors. The Stratified Graph Indexing (SGI) addresses the dataset’s sparsity by centering layers close to the dataset’s center of mass. The index is constructed in each layer as a bi-directional graph and designed to achieve skip list properties, as presented in detail in Section 3. Section 4 offers the novel search and retrieval approach. 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. Note that the SG connections between layers help maintain logarithmic complexity scaling for faster search speed. We present an in-depth comparison of multiple state-of-the-art methods over 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 average Recall and mean average Precision at any depth of retrieval and fast retrieval times, and index size in Section 5. The research paper is concluded with a summary and next steps in Section 9.

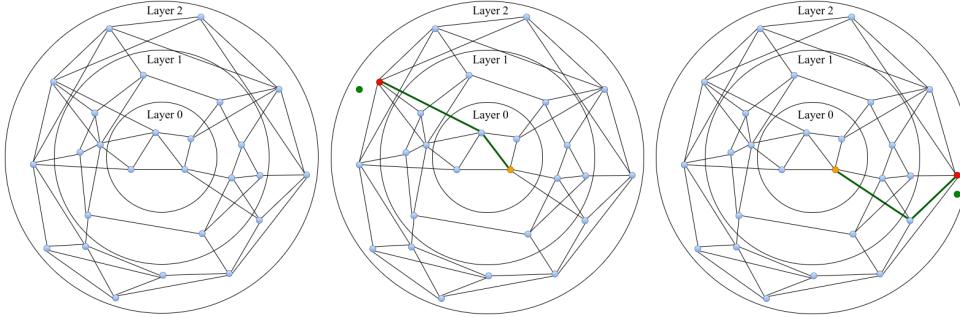


Fig. 5 Illustration of SG index (a) and two different scenarios (b,c) in a Stratified Graph (SG) search. Orange denotes the starting point of a search, red denotes the nearest neighbor, and green edges show the path of the greedy algorithm to the query(shown in green).

3 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 closer to the centroid are stratified in closer layers, while the feature vectors 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 the layers can be illustrated as a set of rotated squares in 2D illustration for the Manhattan distance. The bidirectional graph is constructed by connecting each feature vector to its m -Nearest Neighbors. We specify the maximum number of layers l_{max} as $\log_2 m + 1$. The layer numbers determine the layer width by dividing the distance to the farthest point from the center. Therefore, all the layers have the same width. The index building starts from the outermost layer where all the m neighboring connections for each feature are within the same layer l_{max} . In layer $l_{max} - 1$, each feature vector is connected to $m - 1$ neighbors within layer $l_{max} - 1$ and 1 neighbor with layer l_{max} totaling m neighbors for each feature vectors. Therefore, at any layer l , the number of in-layer connections is determined by $m - (l_{max} - l)$, and next-layer

connections are chosen by $(l_{max} - l)$. 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 leads to a higher index size, slower retrieval, but higher Recall as each point is connected to its actual m nearest neighbors. 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. An outlier filtering factor f filters out elements too far from the mean distance and ensures that outliers do not affect the layer boundaries.

Algorithm 1: LAYER- ING (X, M, f)
<p>Input: data vector P, number of established connections m, outlier filtering factor f</p> <p>Output: element list with assigned layer</p> <pre> 1 $numLayer \leftarrow \lfloor \log_2 m \rfloor$; 2 $cen \leftarrow$ mean of X ; 3 $dist \leftarrow$ distances from the 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 $layeredElem \leftarrow \phi$; 10 foreach (d, p) of $(dist, P)$ do 11 $l \leftarrow \lfloor \frac{d}{r} \rfloor$; 12 add element p to layer l in $layeredElem$; 13 end 14 return $layeredElem$; </pre>

Algorithm 2: SGI ($SGI, P, m, f, cand$)
<p>Input: stratified graph index SGI, dataset P, graph degree m, outlier filtering factor f, size of dynamic candidate list $cand$</p> <p>Output: Updated SGI inserting all the elements in P</p> <pre> 1 $graph \leftarrow \phi$; 2 $layerGraphList \leftarrow \phi$; 3 $layeredElem \leftarrow$ LAYERING(P, m, f) ; 4 foreach $layer$ of $layeredElem$ do 5 $clg \leftarrow \phi$; 6 foreach $elem$ of $layer$ do 7 $clg \leftarrow ADD(clg, elem, m, cand)$; 8 if $layerGraphList$ not empty then foreach g in $layerGraphList$ do $n \leftarrow SGR(g, elem, k = 1, cand)$; 10 $clg \leftarrow$ update clg inserting n to the neighbor list of $elem$; 11 end 12 end 13 end 14 add clg to $layerGraphList$; 15 end 16 $m \leftarrow m - 1$; 17 end 18 $graph \leftarrow$ merge all the graphs in $layerGraphList$; </pre>

Algorithm 1 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, σ . Features with more significant spaces than f times σ from the mean length are filtered out. The remaining elements are then divided into layers based on their distances from the centroid, each containing elements within a specific range of distances. After the layers are determined, bidirectional kNN graphs are constructed for each layer (Algorithm 2 line 4 – 15) and stored in the $layerGraphList$.

Algorithm 2 describes the process of *TODO: missing – Toufik*. The algorithm starts building the bidirectional kNN graphs from the outermost layer. Therefore, all the m neighboring connections are within the same layer for the outermost layer since there is no next layer to the outer layer. For the second to the outermost layer, each feature vector has $m - 1$ neighboring connections within the same layer and 1 neighboring connection with the next-layer feature vector, maintaining a total of m connections for each feature vector. Therefore, as we step towards the innermost layer, the in-layer links decrease (Algorithm 2 line 16), and the next-layer connections increase (Algorithm 2 line 8 – 13), marinating a consistent graph degree of m for all the feature vectors in dataset P . Algorithm 2 line 9-12 ensures the next-layer connections, which helps to capture the global structure of the data. Any new insertion of neighbors in the graph (Algorithm 2 line 10) is simply an SGR search (Algorithm 4) in the existing index. The final graph is constructed by taking a simple union of all the graphs for each layer (Algorithm 2 line 18). The resulting graph captures the local and global structure of the data and can be used for efficient similarity search. The stratified graph indexing (SG) has two phases. During the first phase, iterative insertions add each element simultaneously, 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 has a complexity of $O(|P|\log(|P|))$. Therefore, the overall complexity of the index building of the stratified graph (SG) scales as $O(|P|\log(|P|))$. *TODO: Alg.3 is not referred anywhere in the paper. Why? – Toufik*

4 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(b) and (b), 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(b) and (b)) and top k to the query (marked in green in Figure 5(b) and (b)) 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 question. Algorithm 4 demonstrates how the greedy search process works in the Stratified Graph Retrieval (SGR) technique. Initially, the algorithm locates the

Algorithm 3: ADD ($clg, elem, m, cand$)

Input: current layer graph clg , element to add $elem$, graph degree m , size of dynamic candidate list $cand$

Output: Update clg inserting element $elem$

- 1 $W \leftarrow \text{SGR}(clg, elem, m, cand);$
- 2 $ep \leftarrow \text{get the nearest element from } W \text{ to } elem;$
- 3 **for** $l_c \leftarrow \min(L, l) \dots 0$ **do**
- 4 $W \leftarrow \text{SEARCH-LAYER}(q, ep, \text{efConstruction}, l);$
- 5 neighbors $\leftarrow \text{SELECT-NEIGHBORS}(q, W, M, l_c);$
- 6 **for each** $e \in \text{neighbors}$ **do**
- 7 **if** $|e\text{Conn}| > M_{max}$ **then**
- 8 $e\text{NeuConn} \leftarrow \text{SELECT-NEIGHBORS}(e, \text{Com}, M_{mar}, l_c);$
- 9 $ep \leftarrow W;$
- 10 **if** $I > L$ **then**
- 11 | set enter point for hnsw to q ;
- 12 **end**
- 13 **end**
- 14 **end**
- 15 **end**

Algorithm 4: SGR($g, elem, k, cand$)

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 $ep \leftarrow \text{get entry point of } g$
- 2 $p \leftarrow \text{extract nearest neighbor to } q \text{ starting with } ep$
- 3 $C \leftarrow \text{extract } cand \text{ neighbors to } p \text{ from } g$
 $\text{neighbors} \leftarrow \text{top } k \text{ closest from } C \text{ to } q$
return k -NN
neighbors to q

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 are 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 $O(\log(|P|))$.

Table 1 Dataset characteristics

Dataset	Descriptor Type	Feature Extractor	Dim	DB size in GB	# Instances in millions	% of 0s
VisDrone [30]	Video	[5]	1024	6.2	1.51	69
DOTA2.0 [31]	Image	[5]	1024	11.1	2.69	80
DIOR [32]	Image	[5]	1024	5.2	1.27	83
DEEP10M [18]	Image	[18]	96	3.8	10	0
SIFT10M [33]	Image	[33]	128	0.52	10	25
Crawl840B [34]	Text	[34]	300	5.6	2.2	0

5 Proof Of Concept

In this section, we describe how the Stratified Graph (SG) method is compared to six different state-of-the-art methods: Lightweight approximate Nearest-Neighbor library (N2) [35], Non-Metric Space Library (NMSlib) [36], Hierarchical Navigable Small World library (HNSWlib) [21], Facebook AI Similarity Search library (FaissHNSW) [37], Approximate Nearest-Neighbors Oh Yeah (Annoy) [16], and Hybrid Approximate Nearest-Neighbor Indexing and Search (HANNIS) [10].

Datasets

Table 1 summarizes real data used for the proof of concept. **VisDrone** dataset contains 1,515,007 instances of 1024 dimensional object deep feature descriptors extracted from VisDrone video [30]. **DOTA2.0** dataset contains 2,697,873 instances of 1024-dimensional object deep feature descriptors extracted from DOTA2.0 [31]. **DIOR** dataset contains 1,278,863 instances of 1024 dimensional object deep feature descriptors extracted from DIOR [32]. 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 [18] 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 [8] extracted from Caltech-256 41 \prod 41 whole image patches [33]. Crawl840B dataset with 300 dimensions and 2.2 million instances of vector embeddings of common crawl words using GloVe [34]. 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.

Measures

Evaluation of the indexing and search methods have been conducted from the task’s perspective. Our mission is focused on a visual search for unknown objects or class discovery in the petabytes of image and video archives. Therefore, the Mean Average Precision (MAP@ k) and Average Recall (AR@ k) must remain consistent as k increases while keeping average 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 effects of different methods based on the application. Four performance measurement metrics have been used to evaluate the performance of each technique: MAP@ k , AR@ k , Average Retrieval time, and Index size. The number k is the size of the retrieved set, $k \in [5, 10, 20, 50, 100]$.

Precision@ k for query q , $P_{kq} = (|M_q \cap GT_q|) / |M_q|$.

Mean Average Precision (AP _{k}) is calculated for each query q by considering the Precision at each rank position where relevant items are found in depth k : $AP_{kq} = \sum [P_{kq} \times rel_{kq}] / |GT_q|$. Here, rel_{kq} is a binary indicator equal to 1 if the item at rank k is relevant and 0 otherwise. **MAP@ k** is calculated by averaging the AP values across all queries in the set Q : $MAP@k = \frac{\sum_{q \in Q} AP_{kq}}{|Q|}$

Recall@ k (R@ k) for each query q is defined as: $R_{kq} = (|M_q \cap GT_q|) / |GT_q|$.

AR@ k is a measure of the ability of a retrieval system to retrieve all relevant items

from the entire collection across all queries in the set Q . Once we have the R_{kq} values for all queries in the set Q , we calculate the AR@ k as the average of these Recall values: $AR@k = \frac{1}{|Q|} * \sum_{q \in Q} R_{kq}$.

Index size defines the memory cost to save the indexes in the memory. *TODO: make sure all definitions are consistent – Toufik*

Setup

All experiments were carried out on an Ubuntu 20.04.3 server with 11th generation Intel® Core™ i9-11900K @ 3.5GHzX16 CPU with 128GB RAM and NVIDIA GeForce RTX 3070 8GB mem GPU. The Python implementation of the SG library can be found in <https://anonymous.4open.science/r/SG-4644> *TODO: put this in a reference – Toufik.*

6 Improving Retrieval Effectiveness

Table 2 MAP and AR comparison for VisDrone dataset *TODO: for how many queries? add paper citations here. – Toufik*

Dataset	Method	k	MAP@k					k	AR@k				
			5	10	20	50	100		5	10	20	50	100
Vis Drone	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	0.38		0.86	0.79	0.80	0.88	0.91	
	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	0.80	0.79	0.74	0.61	0.36		0.96	0.95	0.96	0.94	0.91	
DOTA 2.0	N2	0.80	0.88	0.83	0.80	0.66		0.86	0.97	0.97	0.97	0.97	
	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	0.70		0.98	0.98	0.98	0.96	0.96	
	SG	1	1	0.98	0.72	0.51		1	1	0.99	0.97	0.95	
DIOR	N2	1	0.85	0.78	0.73	0.71		1	0.97	0.98	0.99	0.99	
	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	1	1	1	0.98	0.90		1	1	1	1	0.99	
	SG	1	1	1	1	0.93		1	1	1	1	0.99	

Stratified Graph approach is compared with six existing methods for 2.7 million DOTA 2.0, 1.3 DIOR, and 1.5 million VisDrone 1024 dimensional vectors created using [5] *TODO: cite publicly released code too – Toufik*. Table 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 MAP and AR. N2 and HNSWlib perform on par with Stratified Graph (SG) and in terms of effectiveness *only* for larger retrieval sets as illustrated in Table 2. FaissHNSW and Annoy performance quickly degrades for $k > 5$, so the methods are unsuitable for deep descriptor database matching. NMSlib consistently has low MAP and AR for all k in Table 2.

Table 3 MAP and AR comparison for non-deep descriptor databases. *TODO: Make this table look like Table 2 above! – Toufik*

Dataset	Method	k	MAP@k					k	AR@k				
			5	10	20	50	100		5	10	20	50	100
DEEP 10M	N2	0.30	0.30	0.30	0.21	0.27	0.62	0.75	0.79	0.84	0.87		
	NMSlib	0.96	0.81	0.60	0.25	0.12	0.98	0.95	0.90	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	0.87		
	SG	0.78	0.72	0.63	0.52	0.41	0.92	0.88	0.87	0.84	0.84		
SIFT 10M	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	0.42	0.22	0.94	0.93	0.93	0.92	0.90		
	SG	0.94	0.83	0.74	0.42	0.27	0.96	0.98	0.96	0.90	0.82		
Crawl 840B	N2	0.61	0.65	0.65	0.65	0.48	0.66	0.78	0.91	0.95	0.97		
	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	0.94	0.92	0.91	0.83	0.76	0.98	0.97	0.98	0.96	0.96		
	SG	0.70	0.64	0.57	0.41	0.31	0.70	0.70	0.70	0.69	0.70		

Next, let's see if the methods behave similarly to different DNNs 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 [18] were compressed and normalized into 96-dimensional vectors using principal component analysis. MAP@ k and AR@ k -retrievals in Table 3 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, which show inconsistent retrieval performance in Table 3. We also observe an interesting behavior of N2 and Annoy in Table 3: the effectiveness of the indexing method is *improving* with larger k . HNSWlib performs moderately, and FaissHNSW consistently performs poorly in Table 3 compared to SG. Our interpretation from Table 3 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 the most robust algorithm for discovering unknown classes in deep descriptor databases.

Table 3 shows the MAP@ k and AR@ k retrieval results for the SIFT10M dataset for SG and six comparing methods. SG offers the dominating performance in MAP@ 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 AR@ k for higher retrieval results in Table 3. N2 and Annoy show an upward trend in AR@ k retrieval with larger values of k in Table 3. FaissHNSW and NMSlib offer consistently lower performance than SG for both MAP@ k and AR@ k retrieval results in Table 3. HNSWlib performs well and

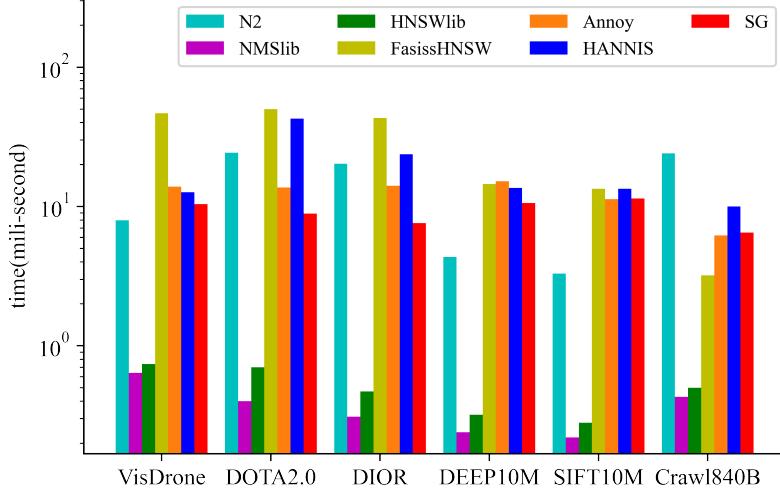


Fig. 6 Average retrieval time for 100 nearest neighbor searches per dataset per method.

achieves similar MAP@ k and AR@ k for higher retrieval results. The MAP@ k and AR@ k retrieval results for the Crawl840B dataset are shown in Table 3. Though our algorithm is specifically designed for deep descriptor database, the performance of SG in Table 3 is compatible with the comparing methods for the vector representation of word embeddings data Crawl840B. In summary, though SG is designed explicitly for similarity search over deep descriptors, its performance is compatible with state-of-the-art algorithms for other descriptor databases.

7 Boosting Indexing Efficiency

Table 4 Index size (GB) per dataset.

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

(SGR) is faster than FaissHNSW, Annoy, and HANNIS library for all six datasets except Crawl840B in Figure 6. 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 (6) 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 extent and type of feature. The Stratified Graph (SG) seems to do well with larger datasets. Table 4 shows the memory cost of saving the indexes in the memory. SG has the most petite index sizes

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 6 shows the promising result on 100 nearest neighbor search. Stratified Graph Retrieval

compared to the comparing methods for all six datasets other than Crawl840B. Here, we use 75 trees for Annoy and 16 neighbor connections for the rest of the six methods. 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 the HNSW algorithm. HNSW arranges the feature vectors in a hierarchical layer of proximity graphs where the upper layers in the 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 most significant 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.

8 Ablation Study

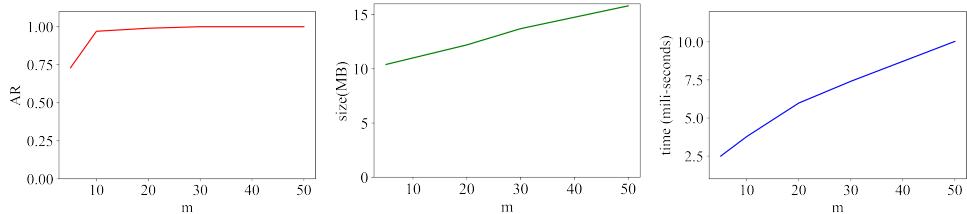


Fig. 7 Demonstration of AR, index size, and retrieval time increasing with a higher degree of m .

The suitable value of neighbor connection m during index building depends on the characteristics of the deep descriptor databases. It can range from 5 to 48, and the ablation study for a random dataset of dimension x and size y is illustrated in Figure 7. The AR increases and approaches 1 as we increase the number of neighbor connections m from 5 to 48 in Figure 7. We achieve an AR of 1 for this random dataset 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 connected m . The lower value of m results in a faster search, losing some accuracy.

9 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 things similar 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 AR and MAP at the cost of slightly higher retrieval times. The MAP and AR 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] Heyse, D.B., Warren, N., Tešić, J.: Identifying maritime vessels at multiple levels of descriptions using deep features. In: Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications, vol. 11006, pp. 423–431 (2019). SPIE
- [2] *Debojyoti Biswas, Jelena Tešić*: Small object difficulty modeling for objects detection in satellite images. In: IEEE 14th International Conference on Computational Intelligence and Communication Networks (CICN), pp. 125–130 (2022). <https://doi.org/10.1109/CICN56167.2022.10008383>
- [3] Zhou, X., Koltun, V., Krähenbühl, P.: Probabilistic two-stage detection. arXiv preprint arXiv:2103.07461 (2021)
- [4] Bochkovskiy, A., Wang, C.-Y., Liao, H.-Y.M.: Yolov4: Optimal speed and accuracy of object detection. arXiv preprint arXiv:2004.10934 (2020)
- [5] *Debojyoti Biswas, Jelena Tešić*: Domain adaptation with contrastive learning for object detection in satellite imagery. IEEE Transactions on Geoscience and Remote Sensing (**under review**) (2023)
- [6] Zhu, X., Lyu, S., Wang, X., Zhao, Q.: 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, pp. 2778–2788 (2021)
- [7] *MMM Rahman, Jelena Tešić*: Hybrid approximate nearest neighbor indexing and search (hannis) for large descriptor databases. In: 2022 IEEE International Conference on Big Data, pp. 3895–3902 (2022). <https://doi.org/10.1109/BigData55660.2022.10020464>
- [8] Lowe, G.: Sift-the scale-invariant feature transform. Int. J **2**(91-110), 2 (2004)
- [9] Arya, S., Mount, D.M., Netanyahu, N.S., Silverman, R., Wu, A.Y.: An optimal algorithm for approximate nearest neighbor searching fixed dimensions. Journal of the ACM (JACM) **45**(6), 891–923 (1998)
- [10] *MMM Rahman, Jelena Tešić*: 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), pp. 6802–6804 (2022). <https://doi.org/10.1109/BigData55660.2022.10021048>
- [11] Fu, C., Wang, C., Cai, D.: High dimensional similarity search with satellite system graph: Efficiency, scalability, and unindexed query compatibility. IEEE Transactions on Pattern Analysis and Machine Intelligence (2021)

- [12] Zhang, J., Ma, R., Song, T., Hua, Y., Xue, Z., Guan, C., Guan, H.: Hierarchical satellite system graph for approximate nearest neighbor search on big data. *ACM/IMS Transactions on Data Science (TDS)* **2**(4), 1–15 (2022)
- [13] Zheng, B., Xi, Z., Weng, L., Hung, N.Q.V., Liu, H., Jensen, C.S.: Pm-lsh: A fast and accurate lsh framework for high-dimensional approximate nn search. *Proceedings of the VLDB Endowment* **13**(5), 643–655 (2020)
- [14] Li, M., Zhang, Y., Sun, Y., Wang, W., Tsang, I.W., Lin, X.: I/o efficient approximate nearest neighbour search based on learned functions. In: 2020 IEEE 36th International Conference on Data Engineering (ICDE), pp. 289–300 (2020). <https://doi.org/10.1109/ICDE48307.2020.00032>
- [15] Kim, S., Yang, H., Kim, M.: Boosted locality sensitive hashing: Discriminative binary codes for source separation. In: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 106–110 (2020). IEEE
- [16] Bernhardsson, E.: Annoy: Approximate Nearest Neighbors in C++/Python. (2018). Python package version 1.17.1. <https://pypi.org/project/annoy/>
- [17] Gallego, A.J., Rico-Juan, J.R., Valero-Mas, J.J.: Efficient k-nearest neighbor search based on clustering and adaptive k values. *Pattern recognition* **122**, 108356 (2022)
- [18] Baranchuk, D., Babenko, A., Malkov, Y.: Revisiting the inverted indices for billion-scale approximate nearest neighbors. *CoRR* **abs/1802.02422** (2018) [1802.02422](https://arxiv.org/abs/1802.02422)
- [19] Jegou, H., Douze, M., Schmid, C.: Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence* **33**(1), 117–128 (2010)
- [20] Li, W., Zhang, Y., Sun, Y., Wang, W., Li, M., Zhang, W., Lin, X.: Approximate nearest neighbor search on high dimensional data: experiments, analyses, and improvement. *IEEE Transactions on Knowledge and Data Engineering* **32**(8), 1475–1488 (2019)
- [21] Malkov, Y.A., Yashunin, D.A.: Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence* **42**(4), 824–836 (2018)
- [22] Fu, C., Xiang, C., Wang, C., Cai, D.: Fast approximate nearest neighbor search with the navigating spreading-out graph. *Proc. VLDB Endow.* **12**(5), 461–474 (2019) <https://doi.org/10.14778/3303753.3303754>
- [23] Beygelzimer, A., Kakade, S., Langford, J.: Cover trees for nearest neighbor. In:

Proceedings of the 23rd International Conference on Machine Learning, pp. 97–104 (2006)

- [24] Fu, C., Cai, D.: Efanna: An extremely fast approximate nearest neighbor search algorithm based on knn graph. arXiv preprint arXiv:1609.07228 (2016)
- [25] Fan, X., Wang, X., Lu, K., Xue, L., Zhao, J.: Tree-based search graph for approximate nearest neighbor search. arXiv preprint arXiv:2201.03237 (2022)
- [26] Iwasaki, M.: Ngt: Neighborhood graph and tree for indexing (2015)
- [27] Iwasaki, M., Miyazaki, D.: Optimization of indexing based on k-nearest neighbor graph for proximity search in high-dimensional data. CoRR **abs/1810.07355** (2018) [1810.07355](https://arxiv.org/abs/1810.07355)
- [28] Iwasaki, M.: Pruned bi-directed k-nearest neighbor graph for proximity search. In: International Conference on Similarity Search and Applications, pp. 20–33 (2016). Springer
- [29] Wang, Y., Ma, H., Wang, D.Z.: Lider: An efficient high-dimensional learned index for large-scale dense passage retrieval. Proc. VLDB Endow. **16**(2), 154–166 (2022) <https://doi.org/10.14778/3565816.3565819>
- [30] Zhu, P., Wen, L., Bian, X., Ling, H., Hu, Q.: Vision meets drones: A challenge. arXiv preprint arXiv:1804.07437 (2018)
- [31] Xia, G.-S., Bai, X., Ding, J., Zhu, Z., Belongie, S., Luo, J., Datcu, M., Pelillo, M., Zhang, L.: Dota: A large-scale dataset for object detection in aerial images. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3974–3983 (2018)
- [32] Li, K., Wan, G., Cheng, G., Meng, L., Han, J.: Object detection in optical remote sensing images: A survey and a new benchmark. ISPRS Journal of Photogrammetry and Remote Sensing **159**, 296–307 (2020)
- [33] Dua, D., Graff, C.: UCI Machine Learning Repository (2017). <http://archive.ics.uci.edu/ml>
- [34] Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 1532–1543 (2014)
- [35] Lee, G.: TOROS N2 - Lightweight Approximate Nearest Neighbor Library Which Runs Fast Even with Large Datasets. (2017). Python package version 0.1.7. <https://github.com/kakao/n2>
- [36] Boytsov, L., Naidan, B.: Engineering efficient and effective non-metric space library. In: International Conference on Similarity Search and Applications, pp.

- 280–293 (2013). Springer
- [37] Johnson, J., Douze, M., Jégou, H.: Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data* **7**(3), 535–547 (2019)