
SPANN: Highly-efficient Billion-scale Approximate Nearest Neighborhood Search

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The in-memory algorithms for approximate nearest neighbor search (ANNS) have
2 achieved great success for fast high-recall search, but are extremely expensive
3 when handling very large scale database. Thus, there is an increasing request for
4 the hybrid ANNS solutions with small memory and inexpensive solid-state drive
5 (SSD). In this paper, we present a simple but efficient memory-disk hybrid indexing
6 and search system, named SPANN, that follows the inverted index methodology. It
7 stores the centroid points of the posting lists in the memory and the large posting
8 lists in the disk. We guarantee both disk-access efficiency (low latency) and
9 high recall by effectively reducing the disk-access number and retrieving high-
10 quality posting lists. In the index-building stage, we adopt a hierarchical balanced
11 clustering algorithm to balance the length of posting lists and augment the posting
12 list by adding the points in the closure of the corresponding clusters. In the search
13 stage, we use a query-aware scheme to dynamically prune the access of unnecessary
14 posting lists. Experiment results demonstrate that SPANN is $2\times$ faster than the
15 state-of-the-art ANNS solution DiskANN to reach the same recall quality 90%
16 with same memory cost in two billion-scale datasets. It can reach 90% recall@1
17 and recall@10 in just around one millisecond with only 32GB memory cost.

18 1 Introduction

19 Vector nearest neighborhood search has played an important role in information retrieval area, such
20 as multimedia search and web search, which provides relevant results by searching vectors with
21 minimum distance to the query vector. Exact solutions for K-nearest neighborhood search [41, 36]
22 are not applicable in big data scenario due to substantial computation cost and high query latency.
23 Therefore, researchers have proposed many kinds of approximate nearest neighborhood search
24 (ANNS) algorithms in the literature [9, 15, 35, 8, 12, 28, 32, 11, 26, 18, 13, 23, 39, 38, 31, 40, 34,
25 29, 16, 24, 7, 10, 21, 42, 17, 33]. However, most of the algorithms mainly focus on how to do low
26 latency and high recall search all in memory with offline pre-built indexes. When targeting to the
27 super large scale vector search scenarios, such as web search, the memory cost will become extremely
28 expensive. There is an increasing request for the hybrid ANNS solutions that use small memory and
29 inexpensive disk to serve the large scale datasets.

30 There are only a few approaches working on the hybrid ANNS solutions, including DiskANN [21]
31 and HM-ANN [33]. Both of them are graph based solutions. DiskANN uses Product Quantization
32 (PQ) [22] to compress the vectors stored in the memory while putting the navigating spread-out
33 graph along with the full-precision vectors on the disk. When a query comes, it traverses the graph
34 according to the distance of quantized vectors and then reranks the candidates according to distance
35 of the full-precision vectors. HM-ANN leverages the heterogeneous memory by placing pivot points
36 in the fast memory and navigable small world graph in the slow memory without data compression.

37 However, it consumes more than 1.5 times larger fast memory than DiskANN. Moreover, the slow
38 memory is still much expensive than disk. Therefore, due to the cheap serving cost, high recall and
39 low latency advantages of DiskANN, it has become the start-of-the-art for indexing billion-scale
40 datasets.

41 In this paper, we argue that the simple inverted index approach can also achieve state-of-the-art
42 performance for large scale datasets in terms of recall, latency and memory cost. We propose SPANN,
43 a simple but surprising efficient memory-disk hybrid vector indexing and search system, that follows
44 the inverted index methodology. SPANN only stores the centroid points of the posting lists in the
45 memory while putting the large posting lists in the disk. We guarantee both low latency and high
46 recall by greatly reducing the number of disk accesses and improving the quality of posting lists.
47 In the index-building stage, we use a hierarchical balanced clustering method to balance the length
48 of posting lists and expand the posting list by adding the points in the closure of the corresponding
49 clusters. In the search stage, we use a query-aware scheme to dynamically prune the access of
50 unnecessary posting lists. Experiment results demonstrate that SPANN is more than two times faster
51 than the state-of-the-art disk-based ANNS algorithm DiskANN to reach the same recall quality 90%
52 with same memory cost in two billion-scale datasets. It can reach 90% recall@1 and recall@10 in
53 just around one millisecond with only 32GB memory cost.

54 2 Background and Related Work

55 Given a set of data vectors $\mathbf{X} \in \mathbb{R}^{n \times m}$ (the data set contains n vectors with m -dimensional features)
56 and a query vector $\mathbf{q} \in \mathbb{R}^m$, the goal of vector search is to find a vector \mathbf{p}^* from \mathbf{X} , called nearest
57 neighbor, such that $\mathbf{p}^* = \arg \min_{\mathbf{p} \in \mathbf{X}} \text{Dist}(\mathbf{p}, \mathbf{q})$. Similarly, we can define K -nearest neighbors.
58 Due to the substantial computation cost and high query latency of the exhaustive search, ANNS
59 algorithms are designed to speedup the search for the approximate K -nearest neighbors in a large
60 dataset in an acceptable amount of time. Most of the ANNS algorithms in the literature mainly focus
61 on the fast high-recall search in the memory. However, with the explosive growth of the vector scale,
62 the memory has become the bottleneck to support large scale vector search. There are only a few
63 approaches working on the ANNS solutions for billion-scale datasets to minimize the memory cost.
64 They can be divided into two categories: inverted index based and graph based methods.

65 The inverted index based methods, such as IVFADC [23], FAISS [24] and IVFOADC+G+P [7],
66 split the vector space into K Voronoi regions by KMeans clustering and only do search in a few
67 regions that are closed to the query. To reduce the memory cost, they use vector quantization, e.g.
68 Product Quantization (PQ) [22], to compress the vectors and store them in the memory. The inverted
69 multi-index (IMI) [5] also uses PQ to compress vectors. It splits the feature space into multiple
70 orthogonal subspaces and constructs a separate codebook for each subspace. The full feature space
71 is produced as a Cartesian product of the corresponding subspaces. Multi-LOPQ[25] uses locally
72 optimized PQ codebook to encode the displacements in the IMI structure. GNO-IMI [6] optimizes
73 the IMI by using non-orthogonal codebooks to produce the centroids. Although they can cut down
74 the memory usage to less than 64GB for one billion 128 dimensional vectors, the recall@1 is very
75 low (only around 60%) due to lossy data compression. Although they can achieve better recall by
76 returning 10 to 100 times more candidates for further reranking, it is often not acceptable in many
77 scenarios.

78 The graph based methods include DiskANN [21] and HM-ANN [33]. Both of them adopt the hybrid
79 solution. DiskANN also stores the PQ compressed vectors in the memory while storing the navigating
80 spread-out graph along with the full-precision vectors on the disk. When a query comes, it traverses
81 the graph using best-first manner according to the distance of quantized vectors and then reranks
82 the candidates according to distance of the full-precision vectors. Similarly, it uses the lossy data
83 compression which will influence the recall quality even though full-precision vector reranking can
84 help retrieve some missing candidates back. The high-cost random disk accesses limit the number of
85 graph traverse and candidate reranking. HM-ANN leverages the heterogeneous memory by placing
86 pivot points promoted by the bottom-up phase in the fast memory and navigable small world graph in
87 the slow memory without data compression. However, it will lead to more than 1.5 times larger fast
88 memory consumption. Moreover, the slow memory is still much expensive than disk and may be not
89 available in some platforms.

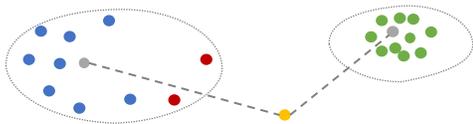


Figure 1: Example of boundary vector missing due to partial search. If we search yellow point, we will search green posting list first since the centroid of green posting list is closer to the yellow point although there are some boundary points (colored red) in the blue posting list that are much closer.

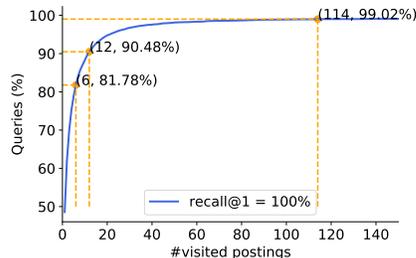


Figure 2: Different queries require different number of posting lists for search. To recall top one result on SIFT1M dataset, we find 80% of queries only need to search 6 posting lists, while 99% of queries need to search 114 posting lists.

90 3 SPANN

91 In this paper, we propose SPANN, a simple but efficient vector indexing and search system, that
 92 follows the inverted index methodology. Different from previous inverted index based methods that
 93 leverage the lossy data compression to reduce the memory cost, SPANN adopts a simple memory-disk
 94 hybrid solution.

95 **Index structure:** The data vectors \mathbf{X} are divided into N posting lists $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$, $\mathbf{X}_1 \cup$
 96 $\mathbf{X}_2 \cup \dots \cup \mathbf{X}_N = \mathbf{X}^1$. The centroids of these posting lists, $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_N$, are stored in the memory
 97 as the fast coarse-grained index that point to the location of the corresponding posting lists in the
 98 disk.

99 **Partial search:** When a query \mathbf{q} comes, we find the K closest centroids, $\{\mathbf{c}_{i_1}, \mathbf{c}_{i_2}, \dots, \mathbf{c}_{i_K}\}$, $K \ll$
 100 N , and load the vectors in the posting lists $\mathbf{X}_{i_1}, \mathbf{X}_{i_2}, \dots, \mathbf{X}_{i_K}$ that correspond to the closest K
 101 centroids into memory for further fine-grained search.

102 3.1 Challenges

103 **Posting length limitation:** Since all the posting lists are stored in the disk, in order to reduce the
 104 disk accesses, we need to bound the length of each posting list so that it can be loaded into memory
 105 in only a few disk reads. This requires us to not only partition the data into a large number of posting
 106 lists but also balance the length of posting lists. This is very difficult due to the substantial high
 107 clustering cost and the balance partition problem itself. The imbalanced posting lists will lead to high
 108 variance of query latency especially when posting lists are stored in the disk.

109 **Boundary issue:** The nearest neighbor vectors of a query \mathbf{q} may locate in the boundary of multiple
 110 posting lists. Since we only search a small number of relevant posting lists, some true neighbors of
 111 \mathbf{q} that located in other posting lists will be missing (Illustrated in Figure 1). If red points are only
 112 represented by the centroid of blue posting list, they will be missing in the nearest neighborhood
 113 search of yellow point.

114 **Diverse search difficulty:** We find that different queries may have different search difficulty. Some
 115 queries only need to be searched in one or two posting lists while some queries require to be searched
 116 in a large number of posting lists (Illustrated in Figure 2). If we search the same number of posting
 117 lists for all queries, it will result in either low recall or long latency.

118 All of the above challenges are the reasons why all of previous inverted index approaches adopt lossy
 119 data compression solution that stores all the compressed vectors and the posting lists in the memory.

120 3.2 Key techniques to address the challenges

121 In this paper, we introduce three key techniques to enable the memory-disk hybrid solution that solve
 122 the above challenges. In the index-building stage, we firstly limit the length of the posting lists to

¹For convenience, we use \mathbf{X} to denote both the matrix and the vector set.

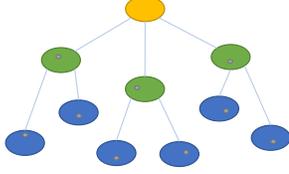


Figure 3: Hierarchical balanced clustering: iteratively balanced partition the vectors in a large cluster (yellow cluster) into a small number of small clusters (green clusters) until each cluster only contains limit number of vectors (blue clusters).



Figure 4: Closure clustering assignment: assign boundary vectors (green points) to multiple closest clusters if its distances to these clusters are nearly the same (blue and yellow clusters).

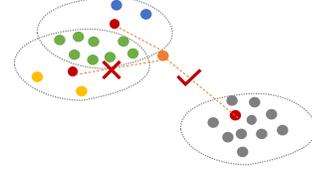


Figure 5: Representative replication: use RNG rule to reduce the similarity of two close posting lists. The orange point will be assigned to blue and grey posting lists although it is closer to yellow list than grey list.

123 effectively reduce the number of disk accesses for each posting list in the online search. Then we
 124 improves the quality of the posting list by expanding the points in the closure of the corresponding
 125 posting lists. This increases the recall probability of the vectors located on the boundary of the
 126 posting lists. In the search stage, we propose a query-aware scheme to dynamically prune the access
 127 of unnecessary posting lists to ensure both high recall and low latency. The detail design of each
 128 technique will be introduced in the following sections.

129 3.2.1 Posting length limitation

130 Limiting the length of posting lists means we need to partition the data vectors \mathbf{X} into a large number
 131 of posting lists $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$. Balancing the length of posting lists means we need to minimize
 132 the variance of the posting length $\sum_{i=1}^N (|\mathbf{X}_i| - |\mathbf{X}|/N)^2$.

133 To address the posting length balance problem, we can leverage the multi-constraint balanced
 134 clustering algorithm [27] to partition the vectors evenly into multiple posting lists:

$$\min_{\mathbf{H}, \mathbf{C}} \|\mathbf{X} - \mathbf{HC}\|_F^2 + \lambda \sum_{i=1}^N \sum_{l=1}^{|\mathbf{X}|} h_{li} - |\mathbf{X}|/N)^2, \text{ s.t. } \sum_{i=1}^N h_{li} = 1. \quad (1)$$

135 where $\mathbf{C} \in \mathbb{R}^{N \times m}$ is the centroids, $\mathbf{H} \in \{0, 1\}^{|\mathbf{X}| \times N}$ represents the cluster assignment, $\sum_{l=1}^{|\mathbf{X}|} h_{li}$
 136 is the number of vectors assigned to the i -th cluster (i.e. $|\mathbf{X}_i|$) and λ is a trade-off hyper parameter
 137 between clustering and balance constraints.

138 However, we find that when the vector number $|\mathbf{X}|$ and the partition number N are very large,
 139 directly using multi-constraint balanced clustering algorithm cannot work due to the difficulty of
 140 large N -partition balanced clustering problem and the extremely high clustering cost. Therefore, we
 141 introduce a hierarchical multi-constraint balanced clustering technique (Figure 3) to not only reduce
 142 the clustering time complexity from $O(|\mathbf{X}| * m * N)$ to $O(|\mathbf{X}| * m * k * \log_k(N))$ (k is a small
 143 constant) but also balance the length of posting lists. We cluster the vectors into a small number
 144 (i.e. k) of clusters iteratively until each posting list contains limit number of vectors. By using this
 145 technique, we can greatly reduce not only the length of each posting list (disk accesses) but also the
 146 index build cost.

147 Moreover, since the number of centroids is very large, finding the nearest posting lists for a query
 148 needs to consume large computation cost. In order to make the navigating computation more
 149 meaningful, we replace the centroid with the vector that is closest to the centroid to represent each
 150 posting list. Then the wasted navigating computation is transformed to the distance computation for a
 151 subset of real candidates.

152 What's more, in order to quickly find a small number of nearest posting lists for a query, we create
 153 a memory SPTAG [10] (MIT license) index for all the vectors that represent the centroids of the
 154 posting lists. SPTAG constructs space partition trees and a relative neighborhood graph as the vector
 155 index which can speedup the nearest centroids search to sub-millisecond response time.

156 **3.2.2 Posting list expansion**

157 To deal with boundary issue, we need to increase the visibility for those vectors that are located in
 158 the boundary of the posting lists. One simple way is to assign each vector to multiple close clusters.
 159 However, it will increase the posting size significantly leading to the heavy disk reads. Therefore, we
 160 introduce a closure multi-cluster assignment solution for boundary vectors on the final clustering step
 161 – assign a vector to multiple closest clusters instead of only the closest one if the distance between the
 162 vector and these clusters are nearly the same (Figure 4 gives an example):

$$\begin{aligned} \mathbf{x} \in \mathbf{X}_{ij} &\iff \text{Dist}(\mathbf{x}, \mathbf{c}_{ij}) \leq (1 + \epsilon) \times \text{Dist}(\mathbf{x}, \mathbf{c}_{i1}), \\ &\text{Dist}(\mathbf{x}, \mathbf{c}_{i1}) \leq \text{Dist}(\mathbf{x}, \mathbf{c}_{i2}) \leq \dots \leq \text{Dist}(\mathbf{x}, \mathbf{c}_{iK}) \end{aligned} \quad (2)$$

163 This means we only duplicate the boundary vectors. For those vectors which are very close to the
 164 centroid of a cluster, they still keep one copy. By doing so, we can effectively limit the capacity
 165 increase due to closure cluster assignment while increasing the recall probability of these boundary
 166 vectors: they will be recalled if any of their closest posting lists is searched.

167 Since each posting list is small and we use closure assignment which will result in some posting lists
 168 that are very close to each other contain the same duplicated vectors (For example, the green vectors
 169 belong to both yellow and blue clusters). Too many duplicated vectors in the close posting lists will
 170 also waste the high-cost disk reads. Therefore, we further optimize the closure clustering assignment
 171 by using RNG rule [37] to choose multiple representative clusters for the assignment of an boundary
 172 vector in order to reduce the similarity of two close posting lists (Figure 5). RNG rule can be simply
 173 defined as we will skip the cluster ij for vector \mathbf{x} if $\text{Dist}(\mathbf{c}_{ij}, \mathbf{x}) > \text{Dist}(\mathbf{c}_{ij-1}, \mathbf{c}_{ij})$. The insight
 174 is two close posting lists are more likely to be both recalled by the navigating index. Instead of
 175 storing similar vectors in close posting lists, it would be better to store different vectors to increase
 176 the number of seen vectors in the online search. From the vector side, it is better to be represented by
 177 posting lists located in different directions (blue and grey posting lists in the example) than just being
 178 represented by posting lists located in the same direction (blue and yellow posting lists).

179 **3.2.3 Query-aware dynamic pruning**

180 In the index-search stage, to process different queries effectively with different resource budget, we
 181 introduce the query-aware dynamic pruning technique to reduce the number of posting lists to be
 182 searched according to the distance between query and centroids. Instead of searching closest K
 183 posting lists for all queries, we dynamically decide a posting list to be searched only if the distance
 184 between its centroid and query is almost the same as the distance between query and the closest
 185 centroid:

$$\begin{aligned} \mathbf{q} \xrightarrow{\text{search}} \mathbf{X}_{ij} &\iff \text{Dist}(\mathbf{q}, \mathbf{c}_{ij}) \leq (1 + \epsilon) \times \text{Dist}(\mathbf{q}, \mathbf{c}_{i1}), \\ &\text{Dist}(\mathbf{q}, \mathbf{c}_{i1}) \leq \text{Dist}(\mathbf{q}, \mathbf{c}_{i2}) \leq \dots \leq \text{Dist}(\mathbf{q}, \mathbf{c}_{iK}) \end{aligned} \quad (3)$$

186 By further reducing those unnecessary posting lists in the closest K posting lists, we can significantly
 187 reduce the query latency while still preserving the high recall by leveraging the resource more
 188 reasonably and effectively.

189 **4 Experiment**

190 In this section we first present the experimental comparison of SPANN with the current state-of-the-art
 191 ANNS algorithms. Then we conduct the ablation studies to further analyze the contribution of each
 192 technique. Finally, we setup an experiment to demonstrate the scalability of SPANN solution in the
 193 distributed search scenario.

194 **4.1 Experiment setup**

195 We conduct all the experiments on a workstation machine with Ubuntu 16.04.6 LTS, which is
 196 equipped with two Intel Xeon 8171M CPU (2600 MHz frequency and 52 CPU cores), 128GB
 197 memory and 2.6TB SSD organized in RAID-0. The datasets we use are as follows:

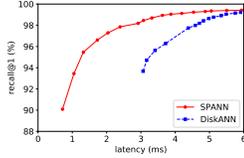


Figure 6: SPANN vs. DiskANN on SIFT1B dataset

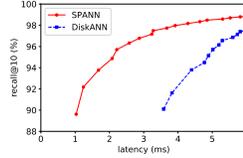


Figure 7: SPANN vs. DiskANN on SPACEV1B dataset

- 198 1. SIFT1M dataset [2] is the most commonly used dataset generated from images for evaluating
 199 the performance of memory-based ANNS algorithms, which contains one million of 128-
 200 dimensional float SIFT descriptors as the base set and 10,00 query descriptors as the test
 201 set.
- 202 2. SIFT1B dataset [2] is a classical dataset for evaluating the performance of ANNS algorithms
 203 that support large scale vector search, which contains one billion of 128-dimensional byte
 204 vectors as the base set and 10,000 query vectors as the test set.
- 205 3. SPACEV1B dataset [4](O-UDA license) is a dataset from commercial search engine which
 206 derives from production data. It represents another different content encoding – deep natural
 207 language encoding. It contains one billion of 100-dimensional byte vectors as a base set and
 208 29,316 query vectors as the test set.

209 The comparison metrics to demonstrate the performance are:

- 210 1. Recall: We compare the R vector ids returned by ANNS with the R ground truth vector ids
 211 to calculate the recall@ R . Since there exist multiple data vectors sharing the same distance
 212 with the query vector, we also replace some of the ground truth vector ids with the vector
 213 ids that sharing the same distance to the query vector in the recall calculation.
- 214 2. Latency: We use the query response time in milliseconds as the query latency.
- 215 3. VQ (Vector-Query): The product of the number of vectors and the number of queries
 216 per second a machine can serve (which is introduced in GRIP [42]). It demonstrates the
 217 serving capacity of the search engine which takes both query latency and memory cost into
 218 consideration. The higher VQ the system has, the less resource cost it consumes. Here we
 219 use the number of vectors per KB \times the number of queries per second as the VQ capacity.

220 4.2 SPANN on single machine

221 In this section, we demonstrate that inverted index based SPANN solution can also achieve the
 222 state-of-the-art performance in terms of recall, latency and memory cost. We first compare SPANN
 223 with the state-of-the-art billion-scale disk-based ANNS algorithms on two billion-scale datasets. Then
 224 we conduct an experiment on SIFT1M dataset to compare the VQ capacity with the start-of-the-art
 225 all-in-memory ANNS algorithms. For all the experiments in this section, we use the following hyper-
 226 parameters for SPANN: 1) use at most 8 closure replicas for each vector in the closure clustering
 227 assignment; 2) limit the max posting list size to 12KB for byte vectors and 48KB for float vectors.
 228 We increase the maximum number of posting lists to be searched in order to get the different recall
 229 quality.

230 4.2.1 Comparison with state-of-the-art billion-scale disk-based ANNS algorithms

231 We choose the state-of-the-art disk-based ANNS algorithms that can support billion-scale datasets as
 232 our comparison targets. we do not compare with HM-ANN [33] since it is not open sourced and the
 233 required PMM hardware may not be available in some platforms. Therefore, we compare SPANN
 234 only with the state-of-the-art billion-scale disk-based ANNS algorithm DiskANN. We use the default
 235 hyper parameters for DiskANN (same as the paper [21] described).

236 We carefully adjust the navigating memory index size of SPANN by choosing suitable number of
 237 posting lists (about 10-12% of total vector number) to ensure both algorithms consume the same
 238 amount of memory (about 32GB). Figure 6 demonstrates the performance for SIFT1B dataset. From

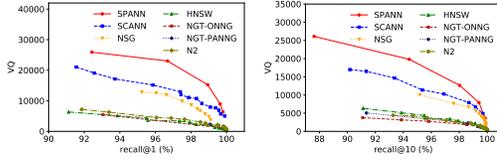


Figure 8: VQ of different ANNS indices

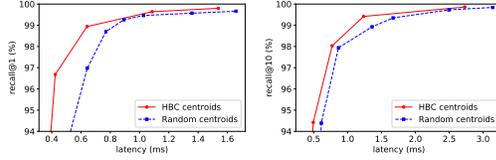


Figure 9: Random VS HBC centroids

239 the results, we find SPANN significantly outperforms DiskANN in both recall@1 and recall@10
 240 especially in the low query latency budget (less than 5ms). Especially, SPANN is more than two
 241 times faster than DiskANN to reach the 95% recall@1 and recall@10.

242 The performance result for SPACEV1B dataset is shown in figure 7. It also demonstrates that SPANN
 243 outperforms DiskANN in both recall@1 and recall@10 when query latency budget is small (less
 244 than 5ms). Especially, DiskANN cannot achieve good recall quality (90%) in less than 4ms, while
 245 SPANN can obtain a recall of 90% in just around 1ms.

246 4.2.2 Comparison with state-of-the-art all-in-memory ANNS algorithms

247 Then we conduct an experiment on SIFT1M dataset to compare the VQ capacity with the start-of-
 248 the-art all-in-memory ANNS algorithms, NSG [16], HNSW [29], SCANN [17], NGT-ONNG [20],
 249 NGT-PANNG [19] and N2 [30]. These algorithms have presented state-of-the-art performance
 250 in the ann-benchmarks [1]. We choose VQ capacity instead of latency as the comparison metric
 251 since these algorithms use much more memory to trade for low latency. However, memory is an
 252 expensive resource which has become the bottleneck for those algorithms to support large scale
 253 datasets. Therefore, we should take both memory and latency into consideration in the performance
 254 comparison.

255 Most of these algorithms are graph based algorithms. For NSG, we get the pre-built index from [3]
 256 and run the performance test with varying SEARCH_L from 1 to 256 which controls the quality of
 257 the search results. For HNSW (nmslib), SCANN, NGT-ONNG, NGT-PANNG and N2 we use the
 258 hyper parameters they provided in the ann-benchmarks [1] that achieve the best performance for the
 259 SIFT1M dataset.

260 Figure 8 demonstrates the VQ capacity of all the algorithms on recall@1 and recall@10. We can see
 261 from the result that SPANN achieves the best VQ capacity consistently across almost all the recall
 262 levels. This means although SPANN cannot achieve as low latency as the all-in-memory ANNS
 263 algorithms due to the high-cost disk accesses during the search, it can obtain the best serving capacity
 264 in the large scale vector search scenario.

265 4.2.3 Ablation studies

266 In this section, we conduct a set of experiments to do the ablation studies on each of our techniques
 267 in the SIFT1M dataset.

268 **Hierarchical balanced clustering** There are two fast ways to partition the vectors on a single
 269 machine into a large number of posting lists: 1) random choose a set of points as the posting list
 270 centroids; 2) using hierarchical balanced clustering (HBC) to generate a set of centroids. We compare
 271 the index quality by generating 16% points as the centroids using these two ways.

272 Figure 9 shows the recall and latency performance of these two centroids choosing algorithms. For
 273 both recall@1 and recall@10, we can see HBC centroid selection is better than random selection.
 274 Moreover, HBC is very fast which clusters one million points into 160K clusters in only around 50
 275 seconds with 64 threads. The whole SPANN index can be built in around 2 minutes.

276 Moreover, how many centroids we need to generate? Small number of centroids can reduce the
 277 navigating memory index size. However, large number of centroids usually means better performance.
 278 Therefore, we need to make reasonable trade-off between the memory usage and the performance.
 279 Figure 10 compares the performance of different number of centroids. From the result, we can see
 280 that the performance will increase significantly with the number of centroids increase when the
 281 number of centroids is small. However, when the number of centroids becomes large enough (16%),

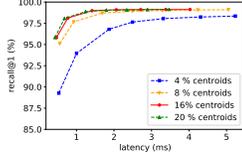


Figure 10: Different number of centroids

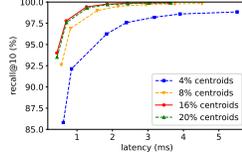


Figure 11: Different number of closure replicas

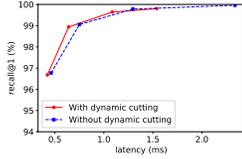
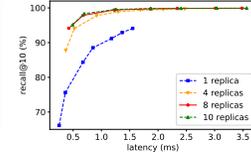
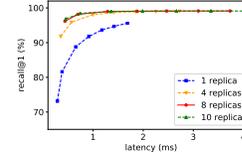


Figure 12: With and without query-aware dynamic pruning

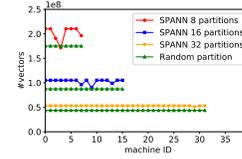
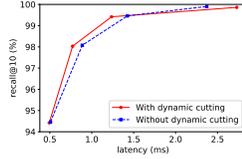
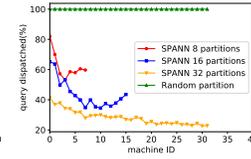


Figure 13: Data size and query access distribution across different machines



282 the performance will not increase any more. Therefore, we can choose 16% of points as the centroids
 283 to achieve both good search performance and small memory usage.

284 **Closure clustering assignment** To use closure clustering assignment, we need to assign a vector to
 285 multiple closed clusters to increase its recall probability during the search. Then at most how many
 286 closure replicas we need to duplicate for a vector to ensure the performance? Too small replicas
 287 cannot help to retrieve those boundary vectors back. However, too many replicas will increase the
 288 posting size greatly which will also affect the performance. Figure 11 demonstrates the performance
 289 of different number of replicas for closure clustering assignment. From the result, we can see that
 290 using more than one replicas improves the performance significantly. However, when the number
 291 of replicas is larger than 8, the performance cannot be improved any more. Therefore, we choose 8
 292 replicas for all of our experiments.

293 **Query-aware dynamic pruning** In order to process different queries effectively during the online
 294 search, we introduce the query-aware dynamic pruning technique to further reduce the number of
 295 posting lists to be searched by pruning those unnecessary posting lists in the closest K posting lists.
 296 We compare the performance with and without query-aware dynamic pruning in the figure 12. From
 297 the result, we can see that with query-aware dynamic pruning, we can further reduce the query latency
 298 without recall drop especially when the latency budget is small. Note that, this technique can reduce
 299 not only the query latency but also the resource usage for a query.

300 4.3 Extension of SPANN to distributed search scenario

301 Compared to the graph base approaches, the additional advantage for inverted index based SPANN
 302 approach is that the partial search on the nearest posting lists idea can be easily extended to the
 303 distributed search scenario, which can handle super large scale vector search with high efficiency and
 304 low serving cost. To demonstrate the scalability of SPANN in distributed search scenario, we partition
 305 the data vectors \mathbf{X} evenly into M partitions $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_M\}$ by using the multi-constraint
 306 balanced clustering and closure clustering assignment techniques in the distributed index build stage,
 307 where M is the number of machines. In the online search stage, we also adopt the query-aware
 308 dynamic pruning technique to reduce the number of dispatched machines, which effectively limits
 309 the total cpu and IO cost for a query.

310 The only challenge for us is that there may have some hot-spot machines. Therefore, we need to
 311 balance not only the data size but also the query access in each machine to avoid the hot spots. To
 312 address the hot-spot challenge, we partition the vectors into multiple small partitions (larger than
 313 machine number) and then use best-fit bin-packing algorithm [14] to pack the small partitions into
 314 large bins (the number of bins equals to the number of machines) according to the history query
 315 access distribution. By doing so, we can effectively balance not only the data size but also the queries
 316 processed on each machine.

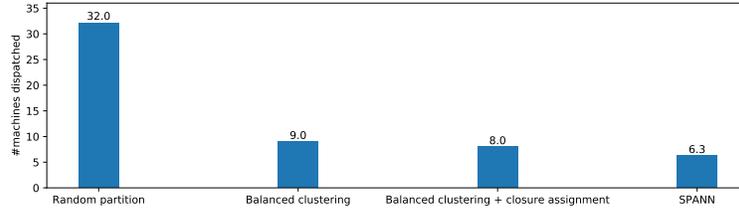


Figure 14: Comparison of #machines dispatched when aggregated recall@10 meets 99%

317 We compare the optimized SPANN solution with traditional random partition and all dispatch solution
 318 to demonstrate the effectiveness of workload reduction and scalability of SPANN in distributed search
 319 scenario. In order to make the result reproducible, we conduct the experiments below based on the
 320 SPACEVIB dataset. In order to focus on the distributed search effectiveness and remove the influence
 321 of ANNS on a single machine, here we use brute-force single machine search results to calculate the
 322 aggregated recall. We use one million queries sampled from the query set as the query access history
 323 and evenly split it into three sets: train, valid and test. The train set is used in offline distributed index
 324 build, and the test set is used in the online search evaluation.

325 4.3.1 Workload reduction and scalability

326 Figure 13 shows the number of vectors and the number of test query accesses in each machine when
 327 partitioning all the base vectors into 8, 16, and 32 partitions. From the result, we can see that SPANN
 328 distributes all the data and query accesses evenly into different machines. Although it increases
 329 the number of vectors in each machine by 20% due to closure assignment, it significantly reduces
 330 the query accesses in each machine compared to the random partition solution. Moreover, SPANN
 331 can continually reduce the query accesses in each machine by using more machines while random
 332 partition cannot. This means we can always add more machines to support more queries per second,
 333 which demonstrates good scalability of our system. The reason why we can achieve good scalability
 334 is that we effectively bound the number of machines to do the search for each query.

335 4.3.2 Analysis

336 Then we analyze how each technique affects the performance. We use 32 partitions case to do the
 337 ablation study. Figure 14 demonstrate the number of machines to search when we require the final
 338 aggregated recall@10 above 99%. For random partition solution, it needs to dispatch the query
 339 to all 32 machines for search. Only use multi-constraints balanced clustering technique, we can
 340 significantly reduce the number of dispatched machines to 9. By adding closure assignment, we can
 341 further reduce the number of dispatched machines to 8. When all the techniques applied (including
 342 query-aware dynamic pruning in the online search), we can finally reduce the number of dispatched
 343 machines to 6.3. This means we can save about 80.3% of computation and IO cost for a query.
 344 Meanwhile, by reducing the number of machines to search for a query, we can further reduce the
 345 query latency since we reduce the number of candidates for final aggregation.

346 5 Conclusion

347 In this paper, we introduce SPANN, a simple but surprising efficient inverted index based ANNS
 348 system, which achieves state-of-the-art performance for large scale datasets in terms of recall,
 349 latency and memory cost. Different from previous inverted index based methods that use lossy data
 350 compression to address the memory bottleneck, SPANN adopts a simple memory-disk hybrid solution
 351 which only stores the centroids of the posting lists in the memory. We guarantee both low latency
 352 and high recall by greatly reducing the number of disk accesses and improving the quality of posting
 353 lists. Experiment results show SPANN can not only establish the new state-of-the-art performance for
 354 billion scale datasets but also achieve good scalability when extended to distributed search scenario.
 355 This demonstrates the ability of hierarchical SPANN to support super large scale vector search with
 356 high efficiency and low serving cost.

References

- 357
- 358 [1] [n.d.]. Benchmarking nearest neighbors. <http://ann-benchmarks.com/>.
- 359 [2] [n.d.]. Datasets for approximate nearest neighbor search. <http://corpus-texmex.irisa.fr/>.
- 360
- 361 [3] [n.d.]. NSG : Navigating Spread-out Graph For Approximate Nearest Neighbor Search. <https://github.com/ZJULearning/nsg>.
- 362
- 363 [4] [n.d.]. SPACEV1B: A billion-Scale vector dataset for text descriptors. <https://github.com/microsoft/SPTAG/tree/master/datasets/SPACEV1B>.
- 364
- 365 [5] Artem Babenko and Victor Lempitsky. 2014. The inverted multi-index. *IEEE transactions on pattern analysis and machine intelligence* 37, 6 (2014), 1247–1260.
- 366
- 367 [6] Artem Babenko and Victor Lempitsky. 2016. Efficient indexing of billion-scale datasets of deep descriptors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2055–2063.
- 368
- 369
- 370 [7] Dmitry Baranchuk, Artem Babenko, and Yury Malkov. 2018. Revisiting the inverted indices for billion-scale approximate nearest neighbors. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 202–216.
- 371
- 372
- 373 [8] J.S. Beis and D.G. Lowe. 1997. Shape indexing using approximate nearest-neighbour search in high-dimensional spaces. In *1997 Conference on Computer Vision and Pattern Recognition {CVPR}'97, June 17-19, 1997, San Juan, Puerto Rico*. 1000–1006.
- 374
- 375
- 376 [9] Jon Louis Bentley. 1975. Multidimensional binary search trees used for associative searching. *Commun. ACM* 18, 9 (1975), 509–517.
- 377
- 378 [10] Qi Chen, Haidong Wang, Mingqin Li, Gang Ren, Scarlett Li, Jeffery Zhu, Jason Li, Chuanjie Liu, Lintao Zhang, and Jingdong Wang. 2018. *SPTAG: A library for fast approximate nearest neighbor search*. <https://github.com/Microsoft/SPTAG>
- 379
- 380
- 381 [11] Sanjoy Dasgupta and Yoav Freund. 2008. Random projection trees and low dimensional manifolds. *Proceedings of the 40th Annual {ACM} Symposium on Theory of Computing, Victoria, British Columbia, Canada, May 17-20, 2008* (2008), 537–546. <https://doi.org/10.1145/1374376.1374452>
- 382
- 383
- 384
- 385 [12] Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S. Mirrokni. 2004. Locality-sensitive Hashing Scheme Based on P-stable Distributions. In *Proceedings of the Twentieth Annual Symposium on Computational Geometry* (Brooklyn, New York, USA) (SCG '04). 253–262.
- 386
- 387
- 388 [13] Wei Dong, Moses Charikar, and Kai Li. 2011. Efficient k-nearest neighbor graph construction for generic similarity measures. In *Proceedings of the 20th International Conference on World Wide Web, WWW 2011, Hyderabad, India, March 28 - April 1, 2011*. 577–586. <https://doi.org/10.1145/1963405.1963487>
- 389
- 390
- 391
- 392 [14] György Dósa and Jiří Sgall. 2014. Optimal analysis of Best Fit bin packing. In *International Colloquium on Automata, Languages, and Programming*. 429–441.
- 393
- 394 [15] Jerome H. Freidman, Jon Louis Bentley, and Raphael Ari Finkel. 1977. An Algorithm for Finding Best Matches in Logarithmic Expected Time. *ACM Trans. Math. Software* 3, 3 (1977), 209–226.
- 395
- 396
- 397 [16] Cong Fu, Chao Xiang, Changxu Wang, and Deng Cai. 2019. Fast Approximate Nearest Neighbor Search With The Navigating Spreading-out Graphs. *PVLDB* 12, 5 (2019), 461 – 474.
- 398
- 399 [17] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. 2020. Accelerating Large-Scale Inference with Anisotropic Vector Quantization. In *International Conference on Machine Learning*. <https://arxiv.org/abs/1908.10396>
- 400
- 401

- 402 [18] Kiana Hajebi, Yasin Abbasi-Yadkori, Hossein Shahbazi, and Hong Zhang. 2011. Fast Approximate Nearest-Neighbor Search with k-Nearest Neighbor Graph. In *IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011*. 1312–1317. <https://doi.org/10.5591/978-1-57735-516-8/IJCAI11-222>
- 407 [19] Masajiro Iwasaki. 2016. Pruned bi-directed k-nearest neighbor graph for proximity search. In *International Conference on Similarity Search and Applications*. Springer, 20–33.
- 409 [20] 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).
- 412 [21] Suhas Jayaram Subramanya, Fnu Devvrit, Harsha Vardhan Simhadri, Ravishankar Krishnawamy, and Rohan Kadekodi. 2019. Rand-nsg: Fast accurate billion-point nearest neighbor search on a single node. *Advances in Neural Information Processing Systems* 32 (2019), 13771–13781.
- 415 [22] 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.
- 418 [23] Hervé Jégou, Romain Tavenard, Matthijs Douze, and Laurent Amsaleg. 2011. Searching in one billion vectors: re-rank with source coding. In *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 861–864.
- 421 [24] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-scale similarity search with GPUs. *arXiv preprint arXiv:1702.08734* (2017).
- 423 [25] Yannis Kalantidis and Yannis Avrithis. 2014. Locally optimized product quantization for approximate nearest neighbor search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2321–2328.
- 426 [26] Brian Kulis and Trevor Darrell. 2009. Learning to Hash with Binary Reconstructive Embeddings.. In *NIPS*, Vol. 22. Citeseer, 1042–1050.
- 428 [27] Hongfu Liu, Ziming Huang, Qi Chen, Mingqin Li, Yun Fu, and Lintao Zhang. 2018. Fast Clustering with Flexible Balance Constraints. In *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 743–750.
- 431 [28] Ting Liu, Andrew W Moore, Alexander Gray, and Ke Yang. 2004. An investigation of practical approximate nearest neighbor algorithms. *Advances in Neural Information Processing Systems 17 [Neural Information Processing Systems, {NIPS} 2004, December 13-18, 2004, Vancouver, British Columbia, Canada]* (2004), 825–832. <http://papers.nips.cc/paper/2666-an-investigation-of-practical-approximate-nearest-neighbor-algorithms>
- 436 [29] Yu A Malkov and Dmitry A Yashunin. 2016. Efficient and robust approximate nearest neighbor search using Hierarchical Navigable Small World graphs. *arXiv preprint arXiv:1603.09320* (2016).
- 439 [30] Yu A 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.
- 442 [31] Marius Muja and David G. Lowe. 2014. Scalable Nearest Neighbour Algorithms for High Dimensional Data. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36, 11 (2014), 2227–2240. <https://doi.org/10.1109/TPAMI.2014.2321376>
- 445 [32] David Nistér and Henrik Stewénus. 2006. Scalable recognition with a vocabulary tree. *2006 {IEEE} Computer Society Conference on Computer Vision and Pattern Recognition {(CVPR) 2006}, 17-22 June 2006, New York, NY, {USA} 2* (2006), 2161–2168. <https://doi.org/10.1109/CVPR.2006.264>

- 449 [33] Jie Ren, Minjia Zhang, and Dong Li. 2020. HM-ANN: Efficient Billion-Point Nearest Neighbor
450 Search on Heterogeneous Memory. *Advances in Neural Information Processing Systems* 33
451 (2020).
- 452 [34] Anshumali Shrivastava and Ping Li. 2014. Asymmetric LSH (ALSH) for Sublinear Time
453 Maximum Inner Product Search (MIPS). In *NIPS*.
- 454 [35] Robert F. Sproull. 1991. Refinements to nearest-neighbor searching in k-dimensional trees.
455 *Algorithmica* 6, 1-6 (1991), 579–589.
- 456 [36] Eric Sadit Tellez, Guillermo Ruiz, and Edgar Chavez. 2016. Singleton indexes for nearest
457 neighbor search. *Information Systems* 60 (2016), 50–68.
- 458 [37] Godfried T Toussaint. 1980. The relative neighbourhood graph of a finite planar set. *Pattern
459 recognition* 12, 4 (1980), 261–268.
- 460 [38] Jingdong Wang and Shipeng Li. 2012. Query-driven iterated neighborhood graph search for
461 large scale indexing. In *Proceedings of the 20th ACM international conference on Multimedia*.
462 ACM, 179–188.
- 463 [39] Jing Wang, Jingdong Wang, Gang Zeng, Zhuowen Tu, Rui Gan, and Shipeng Li. 2012. Scalable
464 k-nn graph construction for visual descriptors. In *Computer Vision and Pattern Recognition
465 (CVPR), 2012 IEEE Conference on*. IEEE, 1106–1113.
- 466 [40] Jingdong Wang, Naiyan Wang, You Jia, Jian Li, Gang Zeng, Hongbin Zha, and Xian Sheng Hua.
467 2014. Ternary-projection trees for approximate nearest neighbor search. *IEEE Transactions
468 on Pattern Analysis and Machine Intelligence* 36, 2 (2014), 388–403. [https://doi.org/10.
469 1109/TPAMI.2013.125](https://doi.org/10.1109/TPAMI.2013.125)
- 470 [41] Peter N Yianilos. 1993. Data Structures and Algorithms for Nearest Neighbor Search in General
471 Metric Spaces. *Proceedings of the Fourth Annual {ACM/SIGACT-SIAM} Symposium on Discrete
472 Algorithms, 25-27 January 1993, Austin, Texas.* (1993), 311–321.
- 473 [42] Minjia Zhang and Yuxiong He. 2019. Grip: Multi-store capacity-optimized high-performance
474 nearest neighbor search for vector search engine. In *Proceedings of the 28th ACM International
475 Conference on Information and Knowledge Management*. 1673–1682.

476 Checklist

- 477 1. For all authors...
- 478 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
479 contributions and scope? [Yes]
- 480 (b) Did you describe the limitations of your work? [Yes]
- 481 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 482 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
483 them? [Yes]
- 484 2. If you are including theoretical results...
- 485 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 486 (b) Did you include complete proofs of all theoretical results? [N/A]
- 487 3. If you ran experiments...
- 488 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
489 perimental results (either in the supplemental material or as a URL)? **We include the
490 data and instructions needed for reproducing our experiments. For code, we will open
491 source the code soon after our internal review.**
- 492 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
493 were chosen)? [Yes]
- 494 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
495 ments multiple times)? [No]

- 496 (d) Did you include the total amount of compute and the type of resources used (e.g., type
497 of GPUs, internal cluster, or cloud provider)? [Yes]
- 498 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 499 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 500 (b) Did you mention the license of the assets? [Yes]
- 501 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 502 (d) Did you discuss whether and how consent was obtained from people whose data you're
503 using/curating? [Yes]
- 504 (e) Did you discuss whether the data you are using/curating contains personally identifiable
505 information or offensive content? [N/A]
- 506 5. If you used crowdsourcing or conducted research with human subjects...
- 507 (a) Did you include the full text of instructions given to participants and screenshots, if
508 applicable? [N/A]
- 509 (b) Did you describe any potential participant risks, with links to Institutional Review
510 Board (IRB) approvals, if applicable? [N/A]
- 511 (c) Did you include the estimated hourly wage paid to participants and the total amount
512 spent on participant compensation? [N/A]