



Carnegie Mellon University

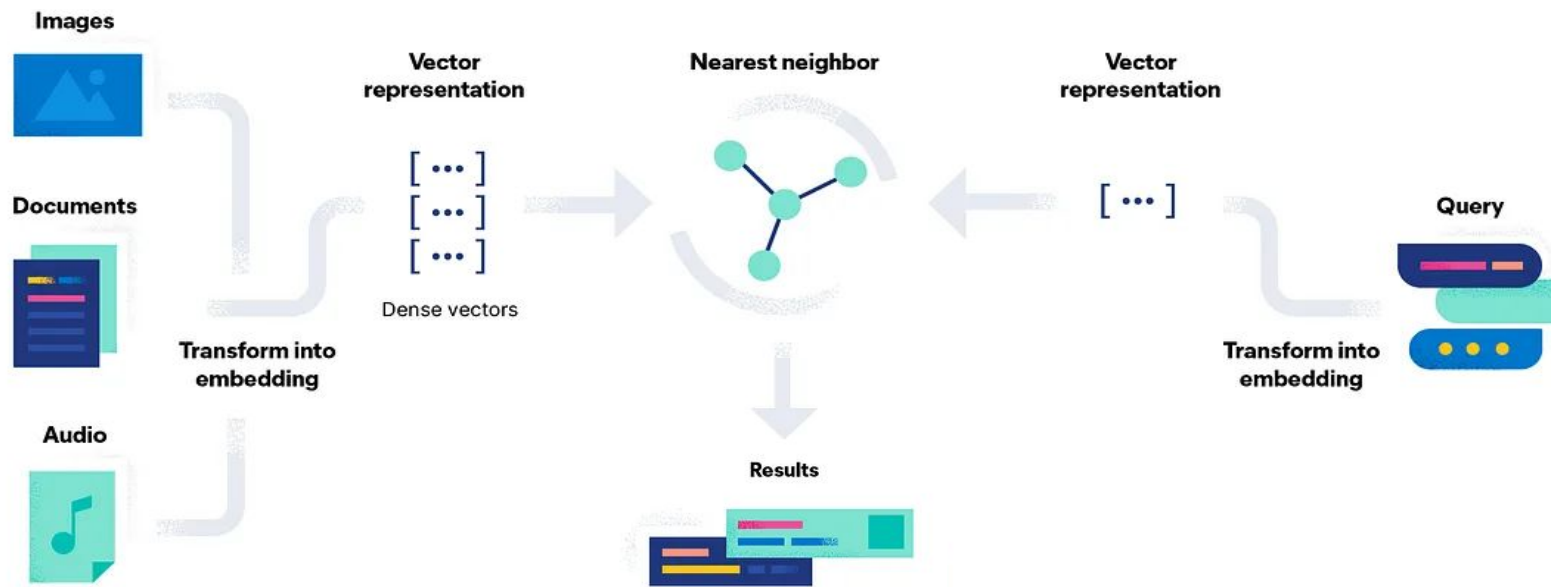
Efficient and robust approximate nearest neighbor search using Hierarchical Navigable Small World (HNSW) graphs

Hu, Tong; Liu, Jiarui; Ma, Christina

4/22/24

Motivation

- Similarity Search: applications in ML, retrieval, and with genAI -> RAG.
- KNN -> ANN: computational complexity vs. search accuracy.

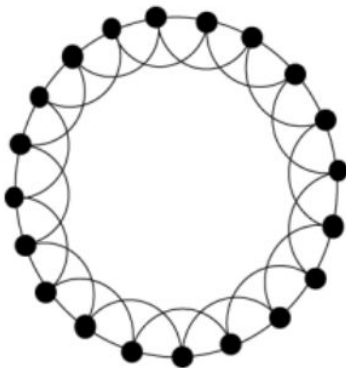


Motivation for Navigable Small Worlds (NSW)

Six degrees of separation experiments run by Milgram in the 1960s.

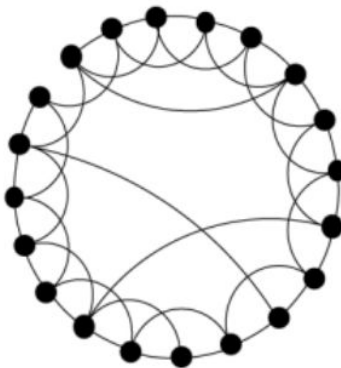


Regular



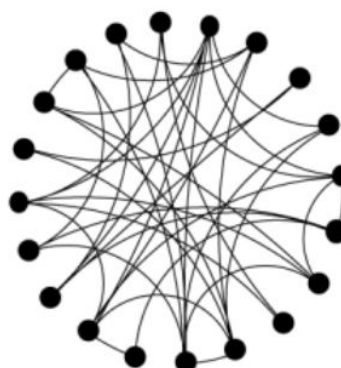
High clustering coefficient
High distance

Small-world



High clustering coefficient
Low distance

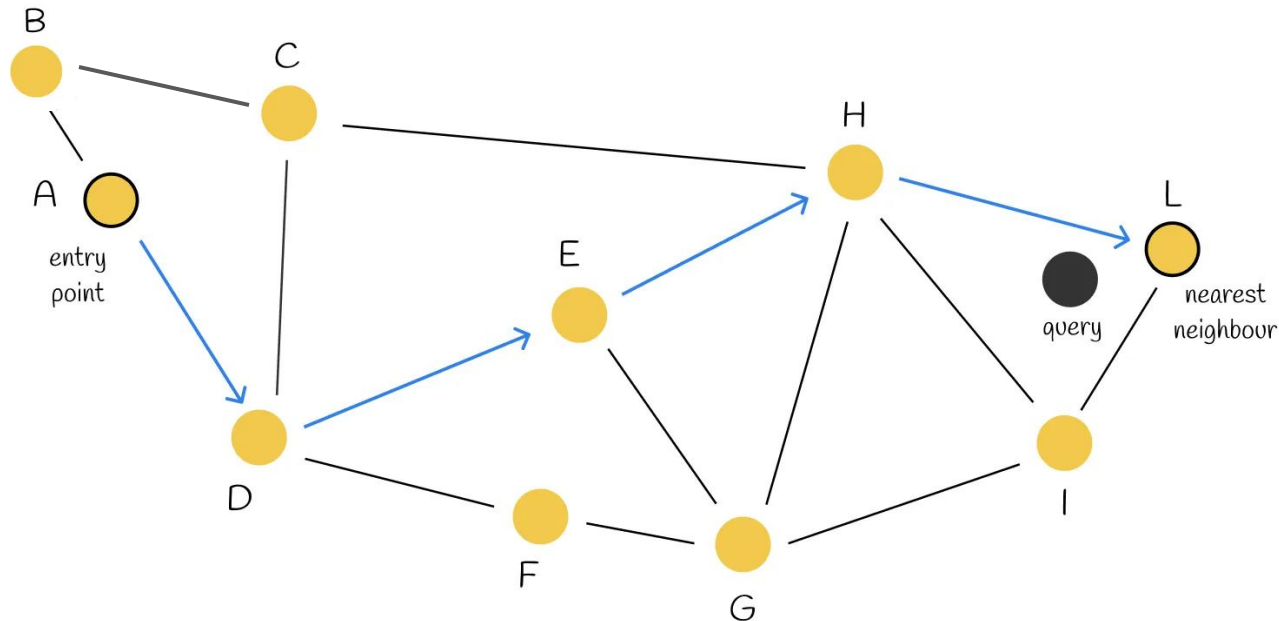
Random



Low clustering coefficient
Low distance

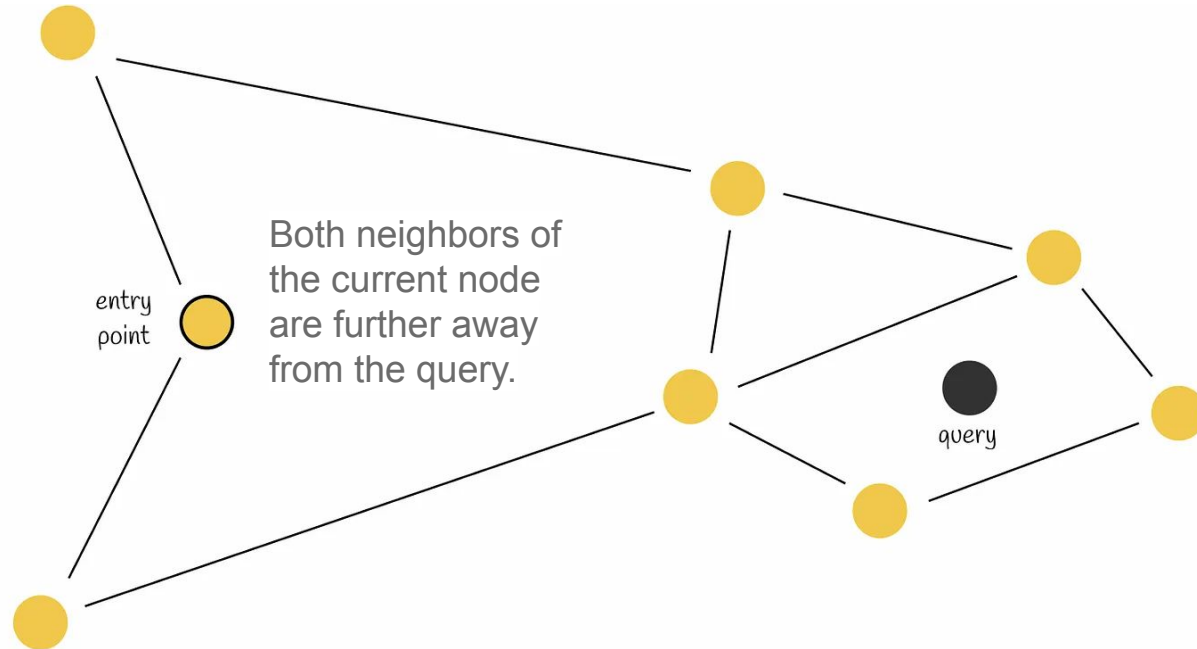
ANN algorithm: Navigable Small Worlds (NSW)

- $O(\log^k n)$ search and insertion, more useful for high dimensional large dataset



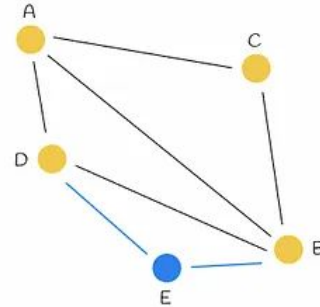
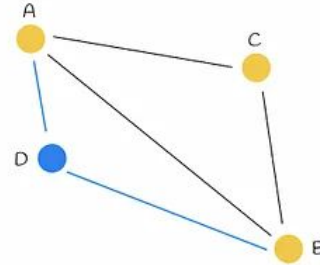
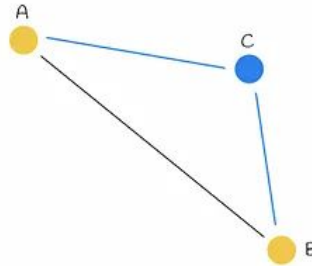
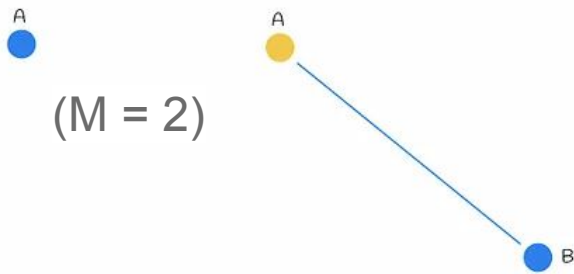
ANN algorithm: Navigable Small Worlds (NSW)

- Greedy search can be trapped in local optimum (early stopping)



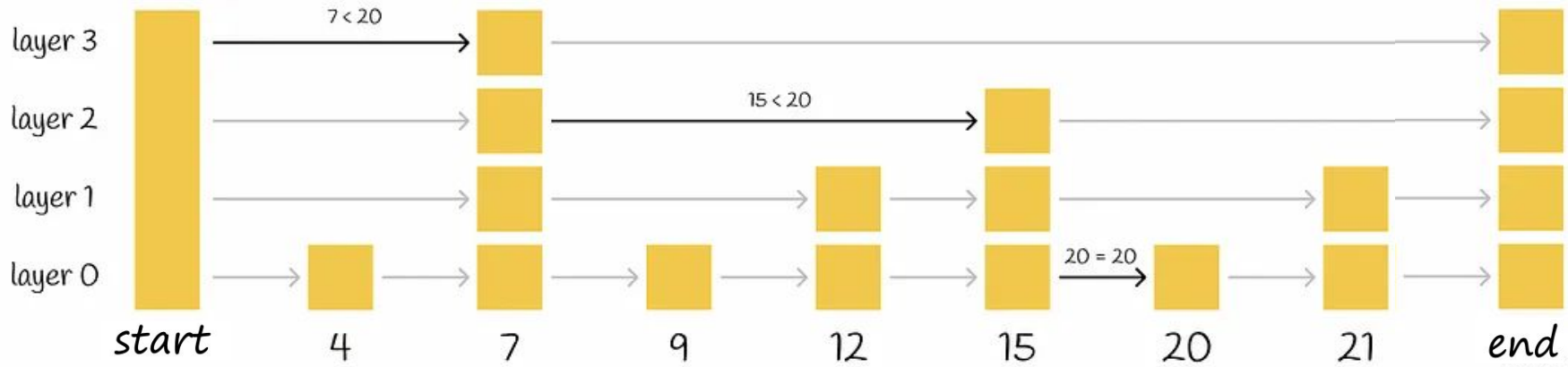
NSW Graph Construction

- Insert random points and link edges to M nearest neighbors (search)
- Longer edges are likely created at the beginning phase of graph construction
 - “later become bridges between the network hubs that keep the overall graph connectivity and allow the logarithmic scaling of the number of hops during greedy routing.”



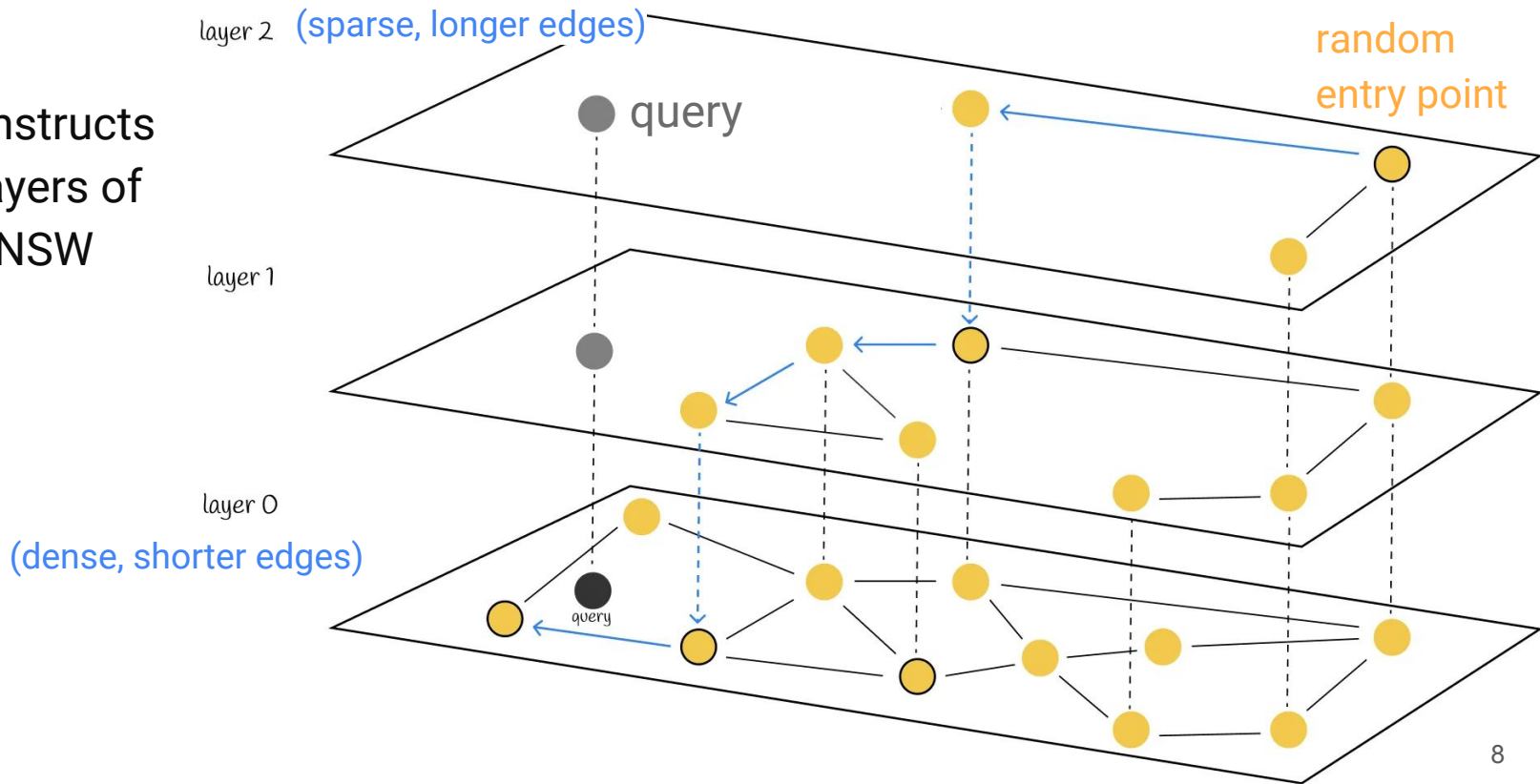
Data structure inspiration: Skip Lists

- $O(\log n)$ time complexity on average for both insertion and search
- Layered format with **longer** edges in the highest layers (for fast search) and **shorter** edges in the lower layers (for accurate search).

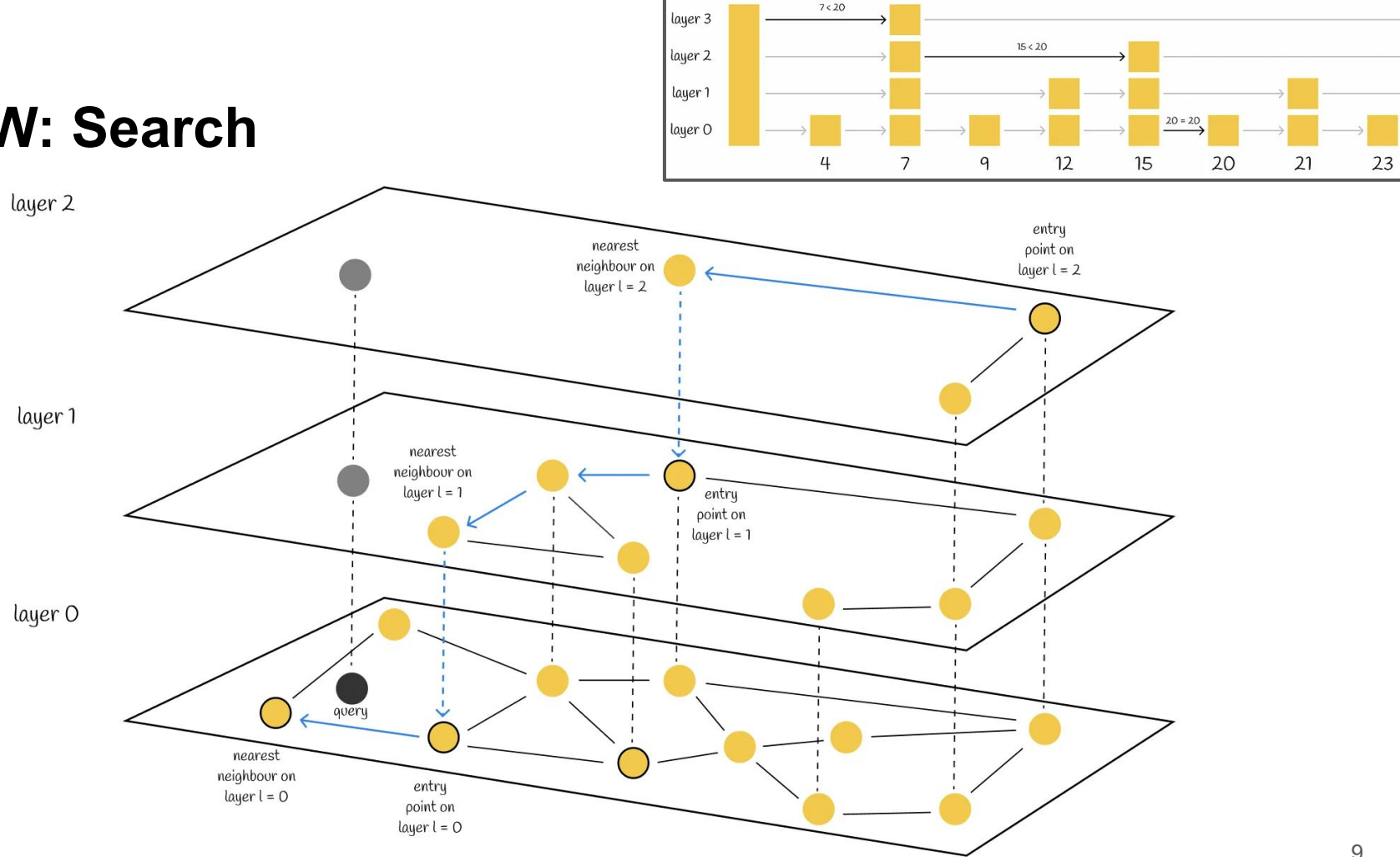


HNSW: *Hierarchical* Navigable Small Worlds

HNSW constructs multiple layers of proximity NSW graphs.



HNSW: Search



HNSW: Search

- Based on principles of skip list and NSW
- Starts from the highest layer
- Proceeds to one level below each time, to find the local nearest neighbor among that layer nodes
- Return the nearest neighbor found on the lowest layer

Algorithm 5

K-NN-SEARCH($hns w, q, K, ef$)

Input: multilayer graph $hns w$, query element q , number of nearest neighbors to return K , size of the dynamic candidate list ef

Output: K nearest elements to q

```
1  $W \leftarrow \emptyset$  // set for the current nearest elements
2  $ep \leftarrow$  get enter point for  $hns w$ 
3  $L \leftarrow$  level of  $ep$  // top layer for  $hns w$ 
4 for  $l_c \leftarrow L \dots 1$ 
5    $W \leftarrow$  SEARCH-LAYER( $q, ep, ef=1, l_c$ )
6    $ep \leftarrow$  get nearest element from  $W$  to  $q$ 
7  $W \leftarrow$  SEARCH-LAYER( $q, ep, ef, l_c=0$ )
8 return  $K$  nearest elements from  $W$  to  $q$ 
```

HNSW: Search

Algorithm 2

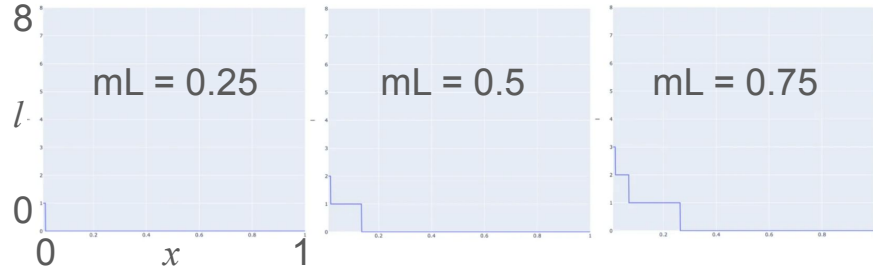
SEARCH-LAYER(q, ep, ef, l_c)

Input: query element q , enter points ep , number of nearest to q elements to return ef , layer number l_c

Output: ef closest neighbors to q

```
1  $v \leftarrow ep$  // set of visited elements
2  $C \leftarrow ep$  // set of candidates
3  $W \leftarrow ep$  // dynamic list of found nearest neighbors
4 while  $|C| > 0$ 
5    $c \leftarrow$  extract nearest element from  $C$  to  $q$ 
6    $f \leftarrow$  get furthest element from  $W$  to  $q$ 
7   if  $distance(c, q) > distance(f, q)$ 
8     break // all elements in  $W$  are evaluated
9   for each  $e \in neighbourhood(c)$  at layer  $l_c$  // update  $C$  and  $W$ 
10    if  $e \notin v$ 
11       $v \leftarrow v \cup e$ 
12     $f \leftarrow$  get furthest element from  $W$  to  $q$ 
13    if  $distance(e, q) < distance(f, q)$  or  $|W| < ef$ 
14       $C \leftarrow C \cup e$ 
15       $W \leftarrow W \cup e$ 
16    if  $|W| > ef$ 
17      remove furthest element from  $W$  to  $q$ 
18 return  $W$ 
```

Graph Construction



Q: How many layers can a node present in the graph?

A: Denote the maximum layer where the node can present as l

$$l = \text{float}[-\ln(\text{uniform}(0, 1)) \cdot mL]$$

The number of layers l for every node is chosen randomly with exponentially decaying probability distribution

“To achieve the optimum performance advantage of the controllable hierarchy, the overlap between neighbors on different layers has to be small.”

mL value tradeoff:

- a smaller mL : more traversals on each layer
- a larger mL : more overlaps

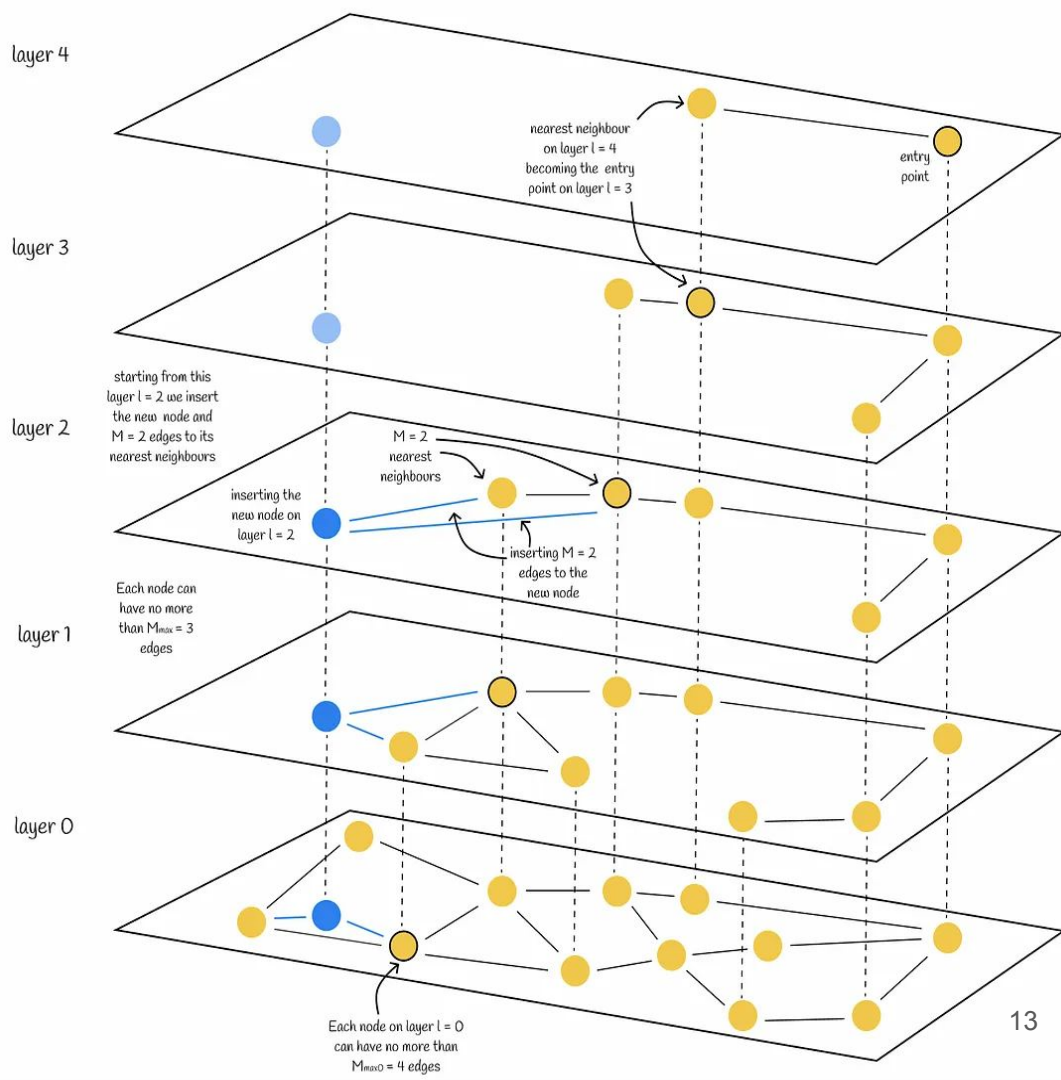
HNSW: Insertion

First phase:

1. Greedily search for the nearest node from the upper layer
2. Use it as an entry point to the next layer until reaching layer l

Second phase:

1. Insert a new node from layer l
2. Greedily search for *efConstruction* nearest neighbors
3. Choose M out of them and build edges
4. Use each of found *efConstruction* nodes as entry points to the next layer until layer 0



HNSW: Insertion

M value tradeoff:

A smaller M is better for lower recalls or low-dimensional data; A larger M is better for high recalls or high-dimensional data

$efConstruction$ value tradeoff:

A larger value implies a more profound search as more candidates are explored, but requires more computations

M_{max} the maximum number of edges a vertex can have

Algorithm 1

INSERT($hns w, q, M, M_{max}, efConstruction, m_L$)

Input: multilayer graph $hns w$, new element q , number of established connections M , maximum number of connections for each element per layer M_{max} , size of the dynamic candidate list $efConstruction$, normalization factor for level generation m_L

Output: update $hns w$ inserting element q

```
1  $W \leftarrow \emptyset$  // list for the currently found nearest elements
2  $ep \leftarrow$  get enter point for  $hns w$ 
3  $L \leftarrow$  level of  $ep$  // top layer for  $hns w$ 
4  $l \leftarrow \lfloor -\ln(\text{unif}(0..1)) \cdot m_L \rfloor$  // new element's level
5 for  $l_c \leftarrow L \dots l+1$ 
6    $W \leftarrow \text{SEARCH-LAYER}(q, ep, ef=1, l_c)$ 
7    $ep \leftarrow$  get the nearest element from  $W$  to  $q$ 
8 for  $l_c \leftarrow \min(L, l) \dots 0$ 
9    $W \leftarrow \text{SEARCH-LAYER}(q, ep, efConstruction, l_c)$ 
10   $neighbors \leftarrow \text{SELECT-NEIGHBORS}(q, W, M, l_c)$  // alg. 3 or alg. 4
11  add bidirectional connections from  $neighbors$  to  $q$  at layer  $l_c$ 
12  for each  $e \in neighbors$  // shrink connections if needed
13     $eConn \leftarrow \text{neighbourhood}(e)$  at layer  $l_c$ 
14    if  $|eConn| > M_{max}$  // shrink connections of  $e$ 
15      // if  $l_c = 0$  then  $M_{max} = M_{max0}$ 
16       $eNewConn \leftarrow \text{SELECT-NEIGHBORS}(e, eConn, M_{max}, l_c)$ 
17      // alg. 3 or alg. 4
18      set  $\text{neighbourhood}(e)$  at layer  $l_c$  to  $eNewConn$ 
19   $ep \leftarrow W$ 
20 if  $l > L$ 
21  set enter point for  $hns w$  to  $q$ 
```

Candidate Selection Simple

Q: Which M nodes to take out of *efConstruction* candidates?

A: Naive way – take M closest candidates

Then a node X is inserted into the graph and needs to be linked to $M=2$ other vertices: B and C.

However, ideally it can be better for navigation if the region A and B can be connected.

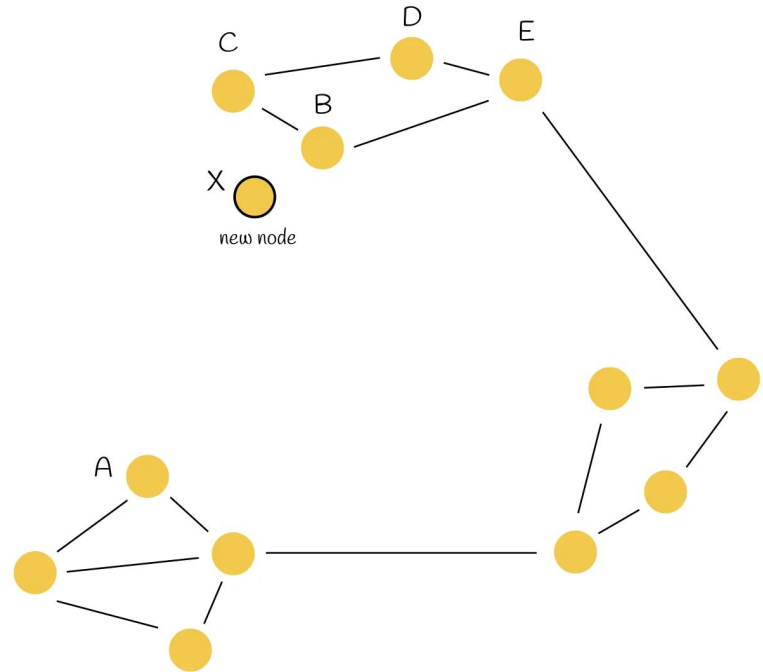
Algorithm 3

SELECT-NEIGHBORS-SIMPLE(q, C, M)

Input: base element q , candidate elements C , number of neighbors to return M

Output: M nearest elements to q

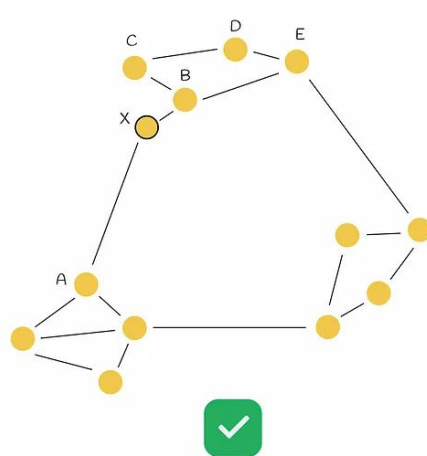
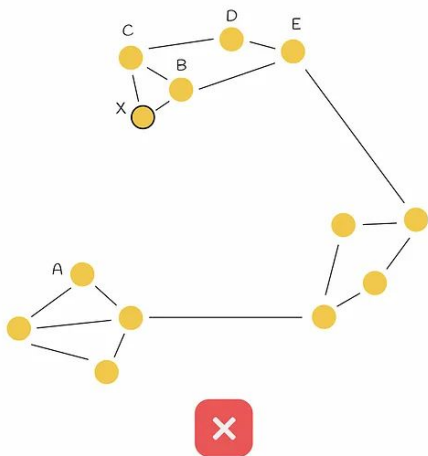
return M nearest elements from C to q



Candidate Selection Heuristic

The heuristic considers both:

- The closest distances between nodes
- The connectivity of different regions on the graph



Algorithm 4

SELECT-NEIGHBORS-HEURISTIC($q, C, M, l_c, \text{extendCandidates}, \text{keepPrunedConnections}$)

Input: base element q , candidate elements C , number of neighbors to return M , layer number l_c , flag indicating whether or not to extend candidate list extendCandidates , flag indicating whether or not to add discarded elements $\text{keepPrunedConnections}$

Output: M elements selected by the heuristic

```
1  $R \leftarrow \emptyset$ 
2  $W \leftarrow C$  // working queue for the candidates
3 if  $\text{extendCandidates}$  // extend candidates by their neighbors
4   for each  $e \in C$ 
5     for each  $e_{adj} \in \text{neighbourhood}(e)$  at layer  $l_c$ 
6       if  $e_{adj} \notin W$ 
7          $W \leftarrow W \cup e_{adj}$ 
8  $W_d \leftarrow \emptyset$  // queue for the discarded candidates
9 while  $|W| > 0$  and  $|R| < M$ 
10   $e \leftarrow$  extract nearest element from  $W$  to  $q$ 
11  if  $e$  is closer to  $q$  compared to any element from  $R$ 
12     $R \leftarrow R \cup e$ 
13  else
14     $W_d \leftarrow W_d \cup e$ 
15  if  $\text{keepPrunedConnections}$  // add some of the discarded
    // connections from  $W_d$ 
16    while  $|W_d| > 0$  and  $|R| < M$ 
17       $R \leftarrow R \cup$  extract nearest element from  $W_d$  to  $q$ 
18 return  $R$ 
```


Complexity Analysis

Search takes $O(\log n)$ time in total

Insertion of a single vertex: $O(\log n)$

HNSW construction requires $O(n * \log n)$ time in total

Evaluation - Implementation

- HNSW implementation uses custom distance functions together with C-style memory management.
- Utilized nmslib implementation of sw-graph for NSW.
- Compare with the most up-to-date SOTA.
- Compare with the SOTA in Euclid Spaces with open-source implementation.

Evaluation - Method

- Comparison with Baseline NSW
- Comparison in Euclid Spaces
- Comparison in General Space
- Comparison with product quantization based algorithms.

Evaluation - HNSW vs. Baseline NSW

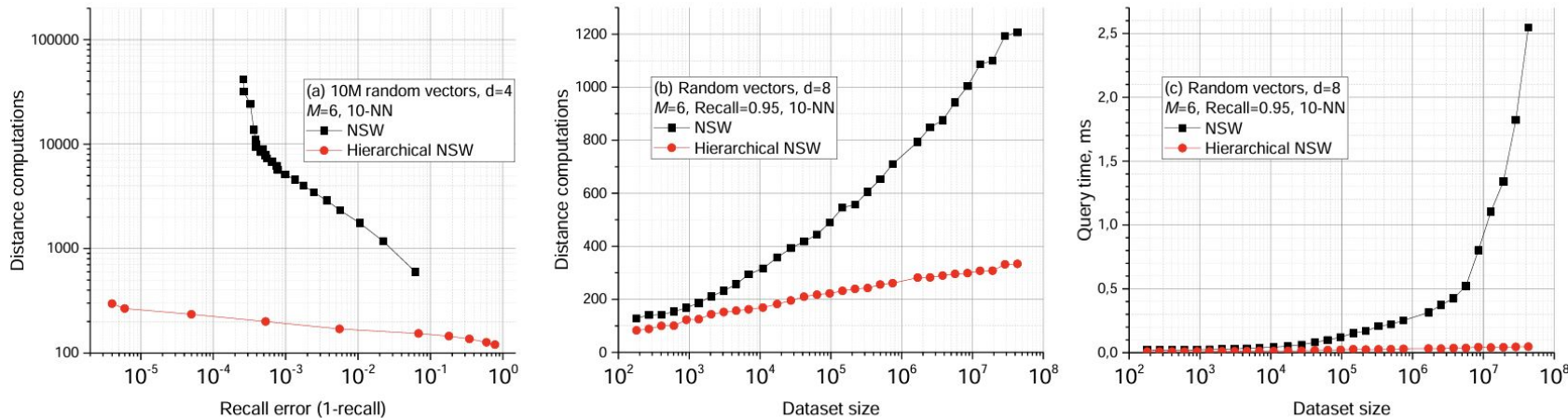


Fig. 12. Comparison between NSW and Hierarchical NSW: (a) distance calculation number vs accuracy tradeoff for a 10 million 4-dimensional random vectors dataset; (b-c) performance scaling in terms of number of distance calculations (b) and raw query(c) time on a 8-dimensional random vectors dataset.

Evaluation - Euclid Spaces - Algorithms to Compare

- Baseline NSW Algorithm
- FLANN
- Annoy
- VP-tree
- FALCONN

Evaluation - Euclid Spaces - Datasets

TABLE 1
Parameters of the used datasets on vector spaces
benchmark.

Dataset	Description	Size	d	BF time	Space
SIFT	Image feature vectors [13]	1M	128	94 ms	L_2
GloVe	Word embeddings trained on tweets [52]	1.2M	100	95 ms	cosine
CoPhIR	MPEG-7 features extracted from the images [53]	2M	272	370 ms	L_2
Random vectors	Random vectors in hypercube	30M	4	590 ms	L_2
DEEP	One million subset of the billion deep image features dataset [14]	1M	96	60 ms	L_2
MNIST	Handwritten digit images [54]	60k	784	22 ms	L_2

Evaluation - Euclid Spaces

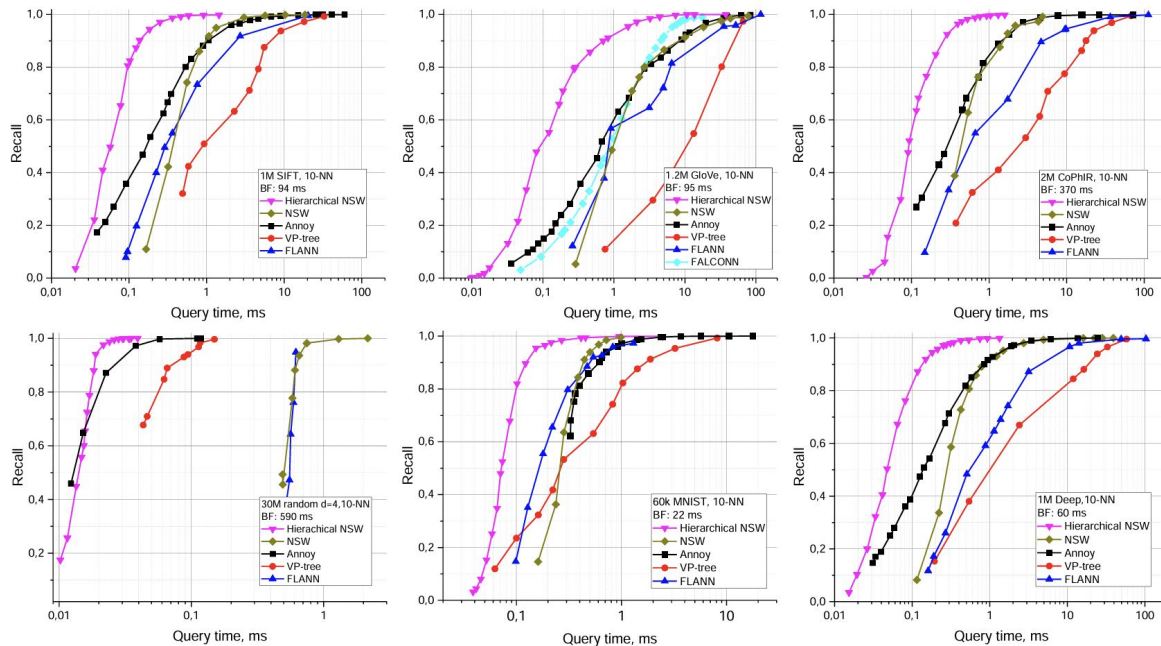


Fig. 13. Results of the comparison of Hierarchical NSW with open source implementations of K-ANNS algorithms on five datasets for 10-NN searches. The time of a brute-force search is denoted as the BF.

Evaluation - General Spaces - Purpose & Algorithms

- Baseline NSW algorithm has several problems on low dimensional datasets as suggested in the paper "Permutation search methods are efficient, yet faster search is possible."
- VP-tree
- Permutation Techniques (NAPP & Brute Force Filtering)
- Baseline NSW Algorithm
- NNDescent-produced proximity graphs

Evaluation - General Spaces - Datasets

TABLE 2.
Used datasets for repetition of the Non-Metric data tests subset.

Dataset	Description	Size	d	BF time	Distance
Wiki-sparse	TF-IDF (term frequency-inverse document frequency) vectors (created via GENSIM [58])	4M	10^5	5.9 s	Sparse cosine
Wiki-8	Topic histograms created from sparse TF-IDF vectors of the wiki-sparse dataset (created via GENSIM [58])	2M	8	-	Jensen-Shannon (JS) divergence
Wiki-128	Topic histograms created from sparse TF-IDF vectors of the wiki-sparse dataset (created via GENSIM [58])	2M	128	1.17 s	Jensen-Shannon (JS) divergence
ImageNet	Signatures extracted from LSVRC-2014 with SQFD (signature quadratic form) distance [59]	1M	272	18.3 s	SQFD
DNA	DNA (deoxyribonucleic acid) dataset sampled from the Human Genome 5 [34].	1M	-	2.4 s	Levenshtein

Evaluation - General Spaces

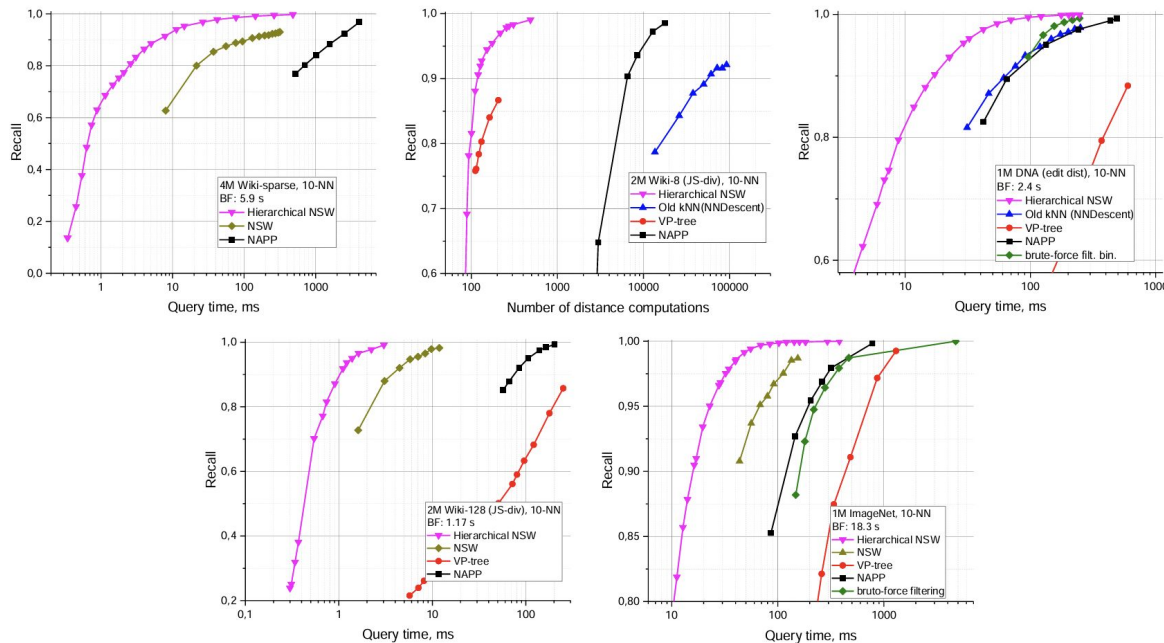


Fig. 14. Results of the comparison of Hierarchical NSW with general space K-ANNS algorithms from the Non Metric Space Library on five datasets for 10-NN searches. The time of a brute-force search is denoted as the BF.

Evaluation - HNSW vs product quantization based algorithms

- PQ-Algorithm: SOTA on billion scale datasets.
- Compare HNSW with SOTA PQ Algorithm in the library: Faiss.

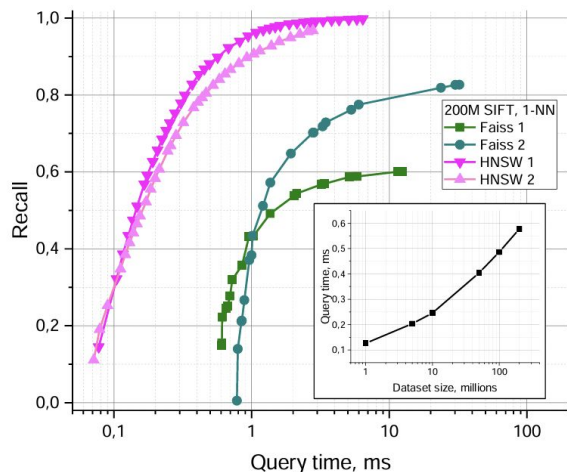


Fig. 15 Results of comparison with Faiss library on the 200M SIFT dataset from [13]. The inset shows the scaling of the query time vs the dataset size for Hierarchical NSW.

TABLE 3.
Parameters for comparison between Hierarchical NSW and Faiss on a 200M subset of 1B SIFT dataset.

Algorithm	Build time	Peak memory (runtime)	Parameters
Hierarchical NSW	5.6 hours	64 Gb	$M=16$, $efConstruction=500$ (1)
Hierarchical NSW	42 minutes	64 Gb	$M=16$, $efConstruction=40$ (2)
Faiss	12 hours	30 Gb	$OPQ64$, $IMI2x14$, $PQ64$ (1)
Faiss	11 hours	23.5 Gb	$OPQ32$, $IMI2x14$, $PQ32$ (2)

Conclusion

- HNSW provides a groundbreaking approach to nearest neighbor search, balancing speed and accuracy effectively even in challenging, high-dimensional spaces.
- The HNSW graph demonstrates robustness to various dataset that was not solvable by baseline NSW. It maintains good performance across different types of datasets without significant tradeoffs.
- This method sets a new benchmark for nearest neighbor searches, offering significant implications for machine learning and data retrieval.

Limitations

- Constructing and maintaining the HNSW graph can consume significant memory, especially for large datasets. This can limit the scalability of the method on memory-constrained systems or for applications with extremely large datasets.
- The search in the HNSW structure always starts from the top layer, thus the structure cannot be easily made distributed like baseline NSW.

Future Work

- The number of added connections per layer M can be a meaningful parameter to tune that strongly affects the construction of the index, thus might improve efficiency and effectiveness of HNSW.
- It would also be interesting to compare HNSW on the full 1B SIFT and 1B DEEP datasets and with functionalities such as element updates and removal.
- Design a distributed pipeline for speedup and memory optimization.

Thanks!