

Towards Optimizing Search and Indexing in Vector Databases

Jayjeet Chakraborty (jayjeetc@ucsc.edu) , Heiner Litz (hlitz@ucsc.edu)



Center for Research
in Systems and Storage



UNIVERSITY OF CALIFORNIA
SANTA CRUZ

Motivation

Vector searches sit on the critical path while querying RAG applications, and hence getting fast responses need vector searches to be highly performant. Also, in hyper-scale setups, vector datasets contain billions of points and require modern hardware for indexing and search operations.

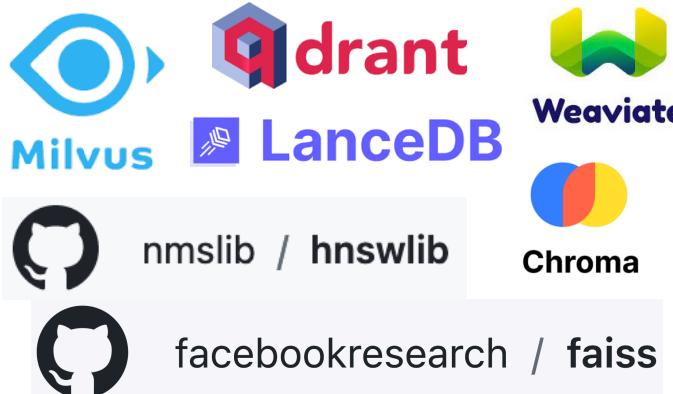
Therefore, having an in-depth understanding of the performance characteristics of vector indexing and search algorithms is crucial.

Vector Search and Indexing

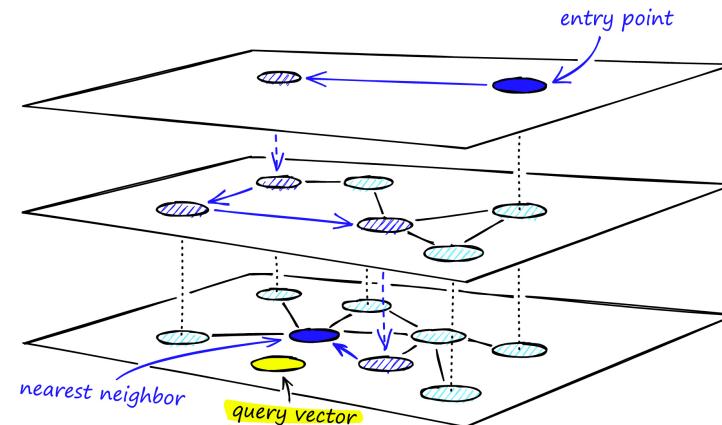
Since, kNN algorithms scale linearly and result in unrealistic search times, ANN algorithms were developed which trade some accuracy for search speed. At the heart of ANN algorithms, are index structures that help narrow down the search space to enable faster searches. Some popular vector indexing algorithms consist of:

- Flat Index
- Inverted File Index (IVF)
- Hierarchical Navigable Small Worlds (HNSW)
- Locality-Sensitive Hashing (LSH)
- Vamana / DiskANN

Vector Databases / Libraries

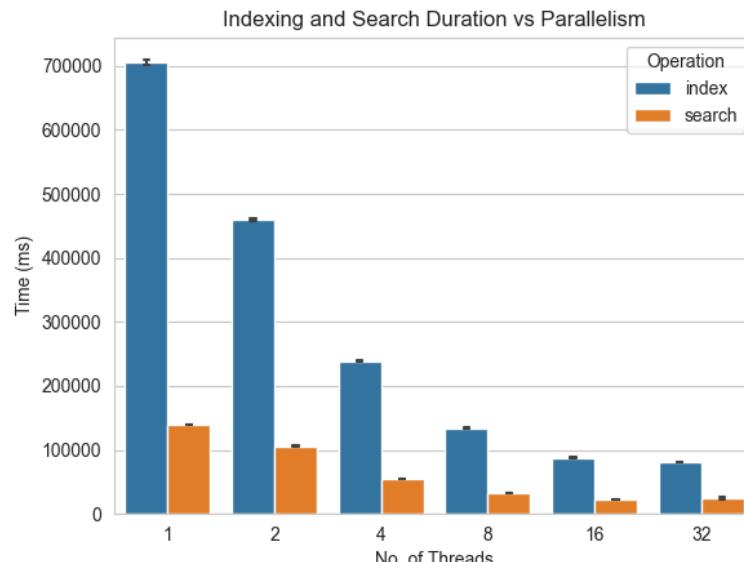


Hierarchical Navigable Small World Graphs

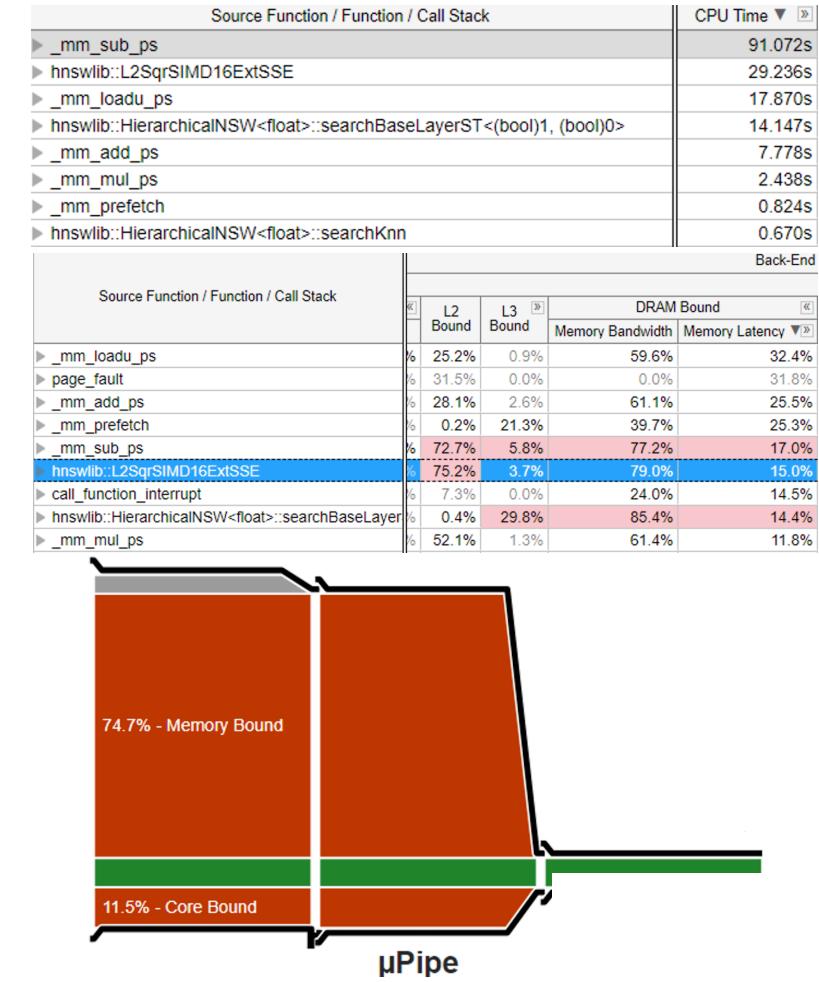


We start from layer 0, pick a entry node, compare its neighbors with the query vector, and move to the neighbor closest to the query vector. Once we find a local minima, we move to that exact node in the next layer and start the search again. The local minima that we find in the last layer is the closest one to our query vector.

Effect of Parallelism in hnswlib



Profiling Vector Searches in hnswlib



Conclusion

Vector search operations can't be batched and can only be parallelized

In-depth understanding of the performance characteristics of vector searches is crucial to find opportunities for acceleration using modern hardware

Vector searches are mostly memory bandwidth bound and most of the query duration is dominated by distance calculations.