LMDiskANN.jl: An Implementation of the Low Memory Disk Approximate Nearest Neighbors Search Algorithm

25 January 2025

Summary

LMDiskANN.jl is a Julia (Bezanson et al. 2017) package that implements the Low Memory Disk Approximate Nearest Neighbor search algorithm (LMDiskANN) (Pan, Sun, and Yu 2023), extending DiskANN-based methods (Jayaram Subramanya et al. 2019; Singh et al. 2021) for fast, accurate billion-point nearest neighbor search while significantly reducing in-memory usage. By leveraging memory-mapped files, a dynamic graph-based index, and tunable BFS expansions, LMDiskANN.jl enables large-scale similarity search on commodity hardware. This package integrates well with embedding-based workflows common in vector databases and modern machine learning pipelines. It also allows for insertions and deletions after the index has been constructed with operations to keep the index pruned of unnecessary connections in the graph.

Features include:

- Low-memory adjacency storage on disk using memory maps.
- Dynamic insert and delete operations on the graph index, adapting to changes in datasets.
- Configurable to tune performance vs. recall.
- Optional user-key LevelDB mapping for flexible ID to and from key lookups.

By combining these capabilities, LMDiskANN.jl aims to reduce the memory footprint and overhead for large-scale nearest neighbor searches. It can be incorporated into any workflow that requires efficient embedding retrieval or similarity search for vectors in high-dimensional spaces.

Statement of need

Approximate Nearest Neighbor (ANN) search is a crucial component in domains such as recommendation systems, information retrieval, and representation learning (e.g., embeddings for natural language or computer vision). Traditional approaches can suffer from excessive memory usage and slow scaling when dealing with billions of points (Nene and Nayar 1997; Wang et al. 2021). By persisting adjacency structures on disk rather than in memory, **LMDiskANN.jl** addresses some of these bottlenecks, providing:

- Reduced Memory Overhead: Only a minimal fraction of data needs to reside in RAM, making it feasible to handle larger datasets on modest machines.
- 2. **Dynamic Updates**: Graph-based insertions and deletions support real-time or streaming scenarios where data is continually changing.
- 3. **High Recall**: Tuning BFS expansions and adjacency degrees can yield high-quality nearest neighbor results.
- 4. **Insertions and Deletions**: Ability to insert and delete from a built index
- 5. **Scalable Architecture**: Built on Julia's high-performance ecosystem, bridging native disk operations and advanced numeric libraries.

This approach benefits practitioners who need large-scale nearest neighbor indexing without specialized cluster infrastructures or extremely large memory capacities. The set up is made to have minimal requirements and has a simple installation procedure. Using the package involves a few number of steps and examples are provided in the documentation.

State of the field

There are various open source ANN implementations and variants in other languages that are standalone or reside within different packages. Within the Julia ecosystem there are fewer options for users to choose from. The package SimilaritySearch.jl (Tellez and Ruiz 2022b, 2022a), offers the most mature codebase. Other notable options exist such as HNSW.jl and NearestNeighborDescent.jl. For applications involving massive datasets that exceed available RAM, these options do not leverage disk space. These alternatives do not offer the user the same flexibility to both insert and delete entries from an already constructed index which is essential for online framework integrations. This was the main motivation for Singh et al. (2021) which describes the need for such utility.

Acknowledgements

Thanks to the Julia community for their continued support of open-source scientific computing. We also acknowledge the authors of LMDiskANN (Pan, Sun, and Yu 2023) and other works it is based on, Jayaram Subramanya et al. (2019) and Singh et al. (2021) for foundational ideas in disk-based graph indexing.

References

Bezanson, Jeff, Alan Edelman, Stefan Karpinski, and Viral B. Shah. 2017. "Julia: A Fresh Approach to Numerical Computing." *SIAM Review* 59 (1): 65–98. https://doi.org/10.1137/141000671.

Jayaram Subramanya, Suhas, Fnu Devvrit, Harsha Vardhan Simhadri, Ravishankar Krishnawamy, and Rohan Kadekodi. 2019. "Diskann: Fast Accurate Billion-Point Nearest Neighbor Search on a Single Node." *Advances in Neural Information Processing Systems* 32.

Nene, Sameer A, and Shree K Nayar. 1997. "A Simple Algorithm for Nearest Neighbor Search in High Dimensions." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19 (9): 989–1003. https://doi.org/10.1109/34.615448.

Pan, Yu, Jianxin Sun, and Hongfeng Yu. 2023. "Lm-Diskann: Low Memory Footprint in Disk-Native Dynamic Graph-Based Ann Indexing." In 2023 Ieee International Conference on Big Data (Bigdata), 5987–96. IEEE. https://doi.org/10.1109/BigData59044.2023.10386517.

Singh, Aditi, Suhas Jayaram Subramanya, Ravishankar Krishnaswamy, and Harsha Vardhan Simhadri. 2021. "Freshdiskann: A Fast and Accurate Graph-Based Ann Index for Streaming Similarity Search." arXiv Preprint arXiv:2105.09613. https://doi.org/10.48550/arXiv.2105.09613.

Tellez, Eric S., and Guillermo Ruiz. 2022a. "SimilaritySearch.jl: Autotuned Nearest Neighbor Indexes for Julia." *Journal of Open Source Software* 7 (75): 4442. https://doi.org/10.21105/joss.04442.

——. 2022b. "Similarity Search on Neighbor Graphs with Automatic Pareto Optimal Performance and Minimum Expected Quality Setups Based on Hyperparameter Optimization." CoRR abs/2201.07917. https://doi.org/10.48550/arXiv.2201.07917.

Wang, Mengzhao, Xiaoliang Xu, Qiang Yue, and Yuxiang Wang. 2021. "A Comprehensive Survey and Experimental Comparison of Graph-Based Approximate Nearest Neighbor Search." arXiv Preprint arXiv:2101.12631. https://doi.org/10.48550/arXiv.2101.12631.