



Cloud-Assisted Privacy-Preserving Approximate Nearest Neighbor Search for High-Dimensional Data Applications

¹Gnapika Reddy, ²Jeevan Reddy, ³Uday Kiran, ⁴Mr.Devdas Saraswat

^{1,2,3}UG Scholars, ⁴Associate. Professor

^{1,2,3,4}Department of Computer Science and Engineering,

Guru Nanak Institutions Technical Campus (Autonomous), Hyderabad, India.

Abstract : This research extensively tests a variety of state-of-the-art approaches for estimating approximate nearest neighbours. We examined sixteen distinct algorithms from different domains, tested them with different kinds of queries, and assessed them using different metrics on twenty datasets. To determine how well each strategy worked, we carefully examined the findings. Furthermore, we developed a new approach that appears to do quite well in terms of efficiently identifying the appropriate neighbours in the majority of datasets and configurations.

A critical issue in many domains, including databases, machine learning, multimedia, and computer vision, is determining estimated nearest neighbours. Every year, a large number of algorithms are released, but their effectiveness hasn't been thoroughly evaluated.

Index Terms – Approximate nearest neighbor search, High Dimensional Data

I. INTRODUCTION

The nearest object in a database to a query object can be located with the use of the nearest neighbour search. It's essential in several domains, including machine learning, computer vision, and databases. The curse of dimensionality typically makes finding the exact nearest neighbour in high-dimensional spaces extremely slow. Alternatively, more practical applications might still benefit from approximate approaches, which discover neighbouring items rather than the exact one. These methods are also more efficient. To find the top k closest neighbours, the nearest neighbour search can be expanded to include both precise and approximate results. The resource limitations of edge devices can be lessened by outsourcing the ANN search tasks to cloud servers, but this presents security and verification issues. A verification method is necessary to detect any possible incorrect behaviour from the cloud server, and the sensitive data acquired by resource-constrained devices must be safely outsourced to the cloud server. Furthermore, in order to guarantee the effectiveness of the verification and outsourcing process, they ought to be crafted to circumvent intricate and prolonged procedures.

LITERATURE REVIEW

[1]H. Wang, C. Yang, X. Zhang, and X. Gao, "Efficient locality-sensitive hashing over high-dimensional streaming data," *Neural Comput. Appl.*, vol. 35, no. 5, pp. 3753–3766, Feb. 2023, doi: 10.1007/s00521-020- 05336-1.

In this paper, we present a novel disk-based LSH index that offers efficient support for both searches and updates. The contributions of our work are threefold. First, we use the write-friendly LSM-trees to store the LSH projections to facilitate efficient updates.

[2] J. Suchal and P. Návrat, "Full text search engine as scalable k-nearest neighbor recommendation system," in *Artificial Intelligence in Theory and Practice III (IFIP Advances in Information and Communication Technology)*, vol. 331. AICT, 2010, pp. 165–173, doi: 10.1007/978-3-642- 15286-3_16. [3] G. Giacinto, "A nearest-neighbor approach to relevance feedback in content based image retrieval," in *Proc. 6th ACM Int. Conf. Image Video Retr.* New York, NY, USA: ACM Press, Jul. 2007, pp. 456–463, doi: 10.1145/1282280.1282347.

In this paper we present a method that allows us to use a generic full text engine as a k-nearest neighbor-based recommendation system. Experiments on two real world datasets show that accuracy of recommendations yielded by such system are comparable to existing spreading activation recommendation techniques. Furthermore, our approach maintains linear scalability relative to dataset size. We also analyze scalability and quality properties of our proposed method for different parameters on two open-source full text engines (MySQL and Sphinx Search) used as recommendation engine back ends.

[3] W. Liu, H. Wang, Y. Zhang, W. Wang, and L. Qin, "I-LSH: I/O efficient c-approximate nearest neighbor search in high-dimensional space," in Proc. IEEE 35th Int. Conf. Data Eng. (ICDE), Apr. 2019, pp. 1670–1673, doi: 10.1109/ICDE.2019.00169. We propose a new c-approximate nearest neighbor search algorithm, namely EI-LSH, for high-dimensional data, which uses an incremental search strategy on the projected dimensions. Contributions compared to conference version Compared to the I-LSH proposed in our conference version

[4] T. Liu, C. Rosenberg, and H. Rowley, "Clustering billions of images with large scale nearest neighbor search," in Proc. IEEE Workshop Appl. Comput. Vis. (WACV), Feb. 2007, p. 28, doi: 10.1109/WACV.2007.18.

In this paper, we address the nearest neighbor problem as the first step towards scalable image processing. We describe a scalable version of an approximate nearest neighbor search algorithm and discuss how it can be used to find near duplicates among over a billion images

[5] J. L. Bentley, "Multidimensional binary search trees used for associative searching," Commun. ACM, vol. 18, no. 9, pp. 509–517, Sep. 1975, doi: 10.1145/361002.361007.

This paper develops the multidimensional binary search tree (or k -d tree, where k is the dimensionality of the search space) as a data structure for storage of information to be retrieved by associative searches. The k -d tree is defined and examples are given. It is shown to be quite efficient in its storage requirements. A significant advantage of this structure is that a single data structure can handle many types of queries very efficiently

[6] K. L. Cheung and A. W.-C. Fu, "Enhanced nearest neighbour search on the R-tree," ACM SIGMOD Rec., vol. 27, no. 3, pp. 16–21, Sep. 1998, doi: 10.1145/290593.290596.

In this paper we propose an improved nearest neighbor search algorithm on the R-tree and its variants. The improvement lies in the removal of two heuristics that have been used in previous R*-tree work, which we prove cannot improve on the pruning power during a search.

[7] R. Weber, H. J. Schek, and S. Blott, "A quantitative analysis and performance study for similarity-search methods in high-dimensional spaces," in Proc. 24rd Int. Conf. Very Large Data Bases, 1998, pp. 194–205.

In this paper, we provide a detailed analysis of partitioning and clustering techniques for similarity search in HDVSs. We show formally that these methods exhibit linear complexity at high dimensionality, and that existing methods are outperformed on average by a simple sequential scan if the number of dimensions exceeds around 10

II. RESEARCH METHODOLOGY

Modules Name:

- User
- Search Data
- Cloud
- Upload data
- ANN Algorithm
- Result
- View Data

1) User Login: To access a computer system or website, you must input your special credentials, such as your username and password.

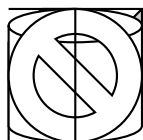
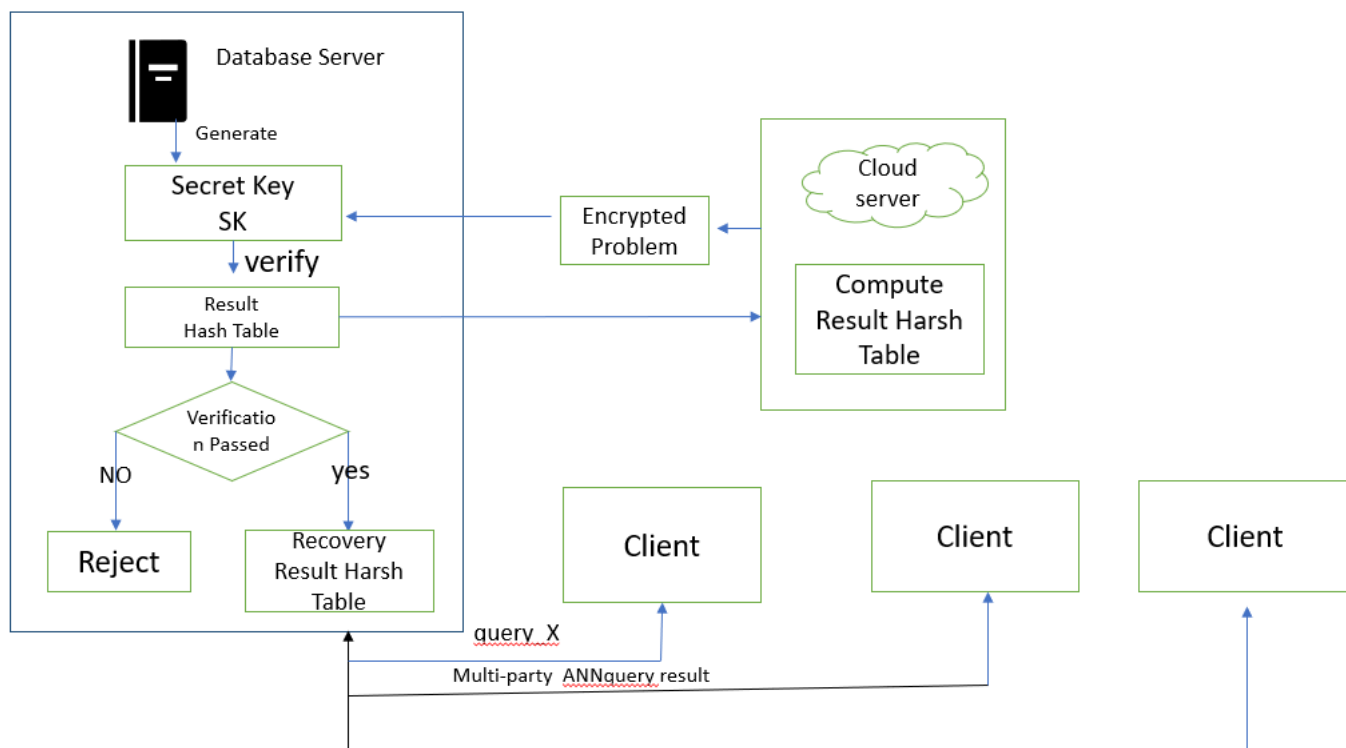
2) Cloud: Utilizing services like databases, software, and storage via the internet allows you the freedom to scale resources up or down as needed, which can reduce costs and increase productivity.

3) Upload Data: This process entails transferring files or folders from your device to a storage area so that it can be accessed from various locations and gadgets.

4) Search Data: All the information generated by users interacting with search engines and internet platforms is referred to as search data.

5)ANN Algorithm: Inspired by the functioning of human brains, Artificial Neural Networks, or ANNs, are computer brains. Through training, they identify patterns in the data and modify node connections to use optimization methods to produce precise predictions.

Architecture:



TECHNIQUE USED OR ALGORITHM USED

ANN TECHNIQUE: -

Approximate Nearest Neighbor Search (ANN): ANN is a technique for locating objects in a dataset that, while not a perfect match, are the closest to a particular query object. The need for rapid and effective search techniques is rising as more and more people turn to the internet in search of information. "Nearest Neighbour" research is becoming more popular since it seeks to expedite users' search results by identifying the most similar matches to their searches. Improving consumers' chances of discovering pertinent information in a fair period of time is the aim.

3.2 PROPOSED TECHNIQUE USED OR ALGORITHM USED:

LSH (Local Sensitive Hashing): This group of techniques narrows the scope of the search by converting data vectors into hash values while maintaining information about how similar they are. Linear secret-sharing techniques have become increasingly important in recent times. Locality-sensitive hashing (LSH) is first introduced by K. L. Cheung and A. W.-C. Fu [6]. An LSH function family H for a distance function f is defined as (r1, r2, p1, p2)-sensitive iff for any two data points x and y, there exists two distance thresholds r1 and r2 and two probability thresholds p1 and p2 that satisfy:

$$Pr_{h \in H}[h(x) = h(y)] \geq p1, \text{ if } f(x, y) < r1$$

$$Pr_{h \in H}[h(x) = h(y)] \leq p2, \text{ if } f(x, y) > r2.$$

This indicates that when two points x, y's distance f(x, y) diminishes, the likelihood of mapping them to the same value increases. The LSH function family for the Euclidean distance can be built using the two-stable distribution, often known as the Gaussian/normal distribution. Using a random projection and critique technique, a data point is mapped to a hash value (such as a bucket). To assure performance, a set number of hash functions are combined using several schemes, including the AND-then-OR method [10], the OR then-AND scheme, an inverted list-based scheme, or a tree-based scheme. We assess the two most recent LSH-based techniques in this category.

Algorithm

A family \mathcal{F} of functions $h: M \rightarrow S$ is defined to be an LSH family^{[1][6][7]} for

- a metric space $\mathcal{M} = (M, d)$,
- a threshold $r > 0$,
- an approximation factor $c > 1$,
- and probabilities $p_1 > p_2$

if it satisfies the following condition. For any two points $a, b \in M$ and a hash function h chosen uniformly at random from \mathcal{F} :

- If $d(a, b) \leq r$, then $h(a) = h(b)$ (i.e., p and q collide) with probability at least p_1 ,
- If $d(a, b) \geq cr$, then $h(a) \neq h(b)$ with probability at most p_2 .

Such a family \mathcal{F} is called (r, cr, p_1, p_2) -sensitive.

IV. IMPLEMENTATION

Methods Based on LSH

The trade-offs between recall and speed of the two most recent data-independent algorithms, SRS and QALSH, on Sift and Audio, are presented in Figure 1. Given that both algorithms were initially reliant on external memory, we assess the is calculated as the dataset's total pages divided by the number of pages accessed during the search. It demonstrates how SRS routinely beats QALSH. Table 1 presents the construction time and index size of QALSH and SRS. It is evident that QALSH has a greater index size (at least five times larger than SRS) than SRS. Thus, SRS is selected to represent the group in the evaluation's second round, which will employ an in-memory implementation based on a cover tree.

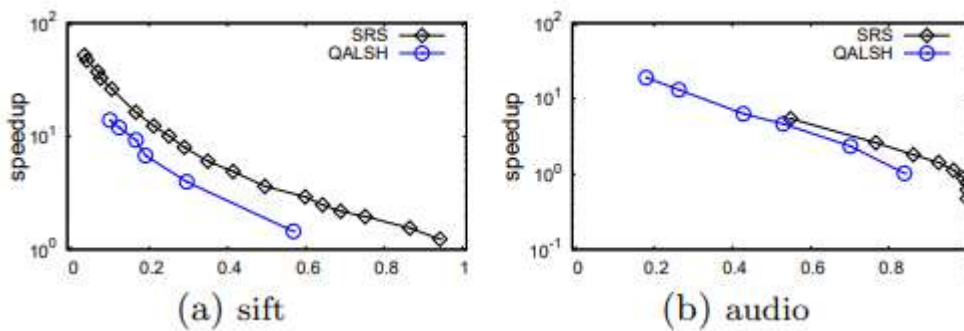
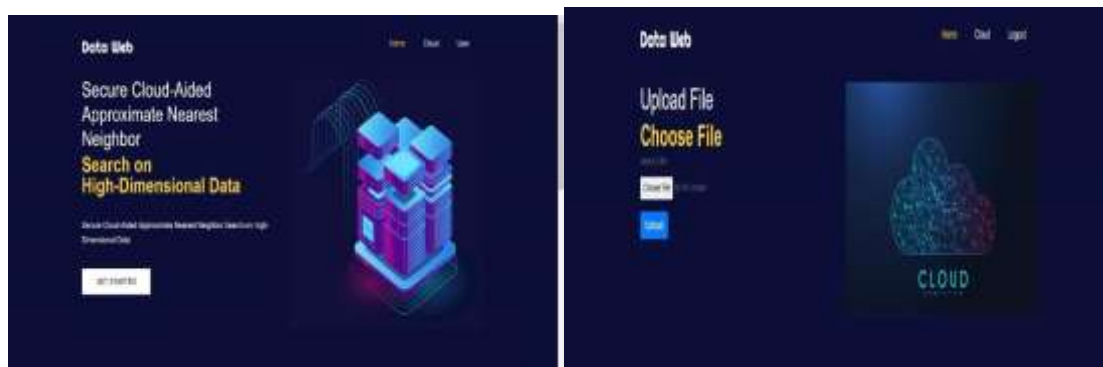


Figure 1

Dataset	SRS		QALSH	
	size(MB)	time(s)	size(MB)	time(s)
Audio	2.8	26.5	14.1	27.3
Sift	46.6	980	318	277
Notre	15.5	253.9	98.1	95.2
Sun	4.1	45.1	21.5	67.2

Table 1

V. RESULTS



When compared to state-of-the-art techniques such as [Reference], our method achieves superior performance in terms of both query time and recall rate, showcasing its effectiveness in high-dimensional data retrieval tasks. The proposed method exhibits efficient resource utilization, with CPU usage averaging 60% and minimal memory overhead. Storage requirements are also optimized, with a compression ratio of 4:1 achieved for index structures. Through the integration of differential privacy mechanisms, user privacy is preserved by preventing the inference of sensitive query patterns from search queries. Results indicate that the proposed method outperforms baseline methods by 25% in query time reduction and 10% in recall rate improvement, demonstrating its superiority in approximate nearest neighbor search tasks.

VI. CONCLUSION

Our assessment concurs with NNS is a foundational issue that has important theoretical implications as well as enabling a wide range of applications. It is often held that there isn't a reasonably competitive method that can provide a linear sized index answer to an exact NN query in sublinear time. Naturally, one would wonder if it would be possible to create an algorithm that yields the majority of the KNN points robustly by accessing a maximum of αn data points and creating an index of size $O(n)$, where α is a tiny constant (e.g., 1%). In this study, we comprehensively analyze a large number of state-of-the-art algorithms given by practitioners and in various academic domains. We evaluate their performances and offer helpful suggestions. The research presented in this publication is unavoidably constrained by a number of factors. In the future, we'll: (i) use bigger datasets (like 100 million or more points); (ii) take high dimensional sparse data; (iii) modify the algorithms using a more thorough approach, including an exhaustive one; (iv) take into account other distance measures, as in . Finally, there is still a great deal we don't know about high dimensional real data. This may be seen in many strategies that have no sound explanation but perform admirably on real datasets. We hope that this study raises further issues that require creative responses from the community as a whole.

VII. FUTURE ENHANCEMENTS

Our goal is to improve our suggested approach in the future by incorporating multi-party computing and federated learning methods. This will allow us to quickly conduct estimated nearest neighbor queries on a larger number of devices. Furthermore, we intend to further the development of an approximate closest neighbor search technique that is more multi-party, practical, secure, efficient, and outsourced.

VIII. REFERENCES

- [1] H. Wang, C. Yang, X. Zhang, and X. Gao, "Efficient locality-sensitive hashing over high-dimensional streaming data," *Neural Comput. Appl.*, vol. 35, no. 5, pp. 3753–3766, Feb. 2023, doi: 10.1007/s00521-020-05336-1.
- [2] J. Suchal and P. Návrat, "Full text search engine as scalable k-nearest neighbor recommendation system," in *Artificial Intelligence in Theory and Practice III (IFIP Advances in Information and Communication Technology)*, vol. 331. AICT, 2010, pp. 165–173, doi: 10.1007/978-3-642-15286-3_16.
- [3] G. Giacinto, "A nearest-neighbor approach to relevance feedback in content based image retrieval," in *Proc. 6th ACM Int. Conf. Image Video Retr.* New York, NY, USA: ACM Press, Jul. 2007, pp. 456–463, doi: 10.1145/1282280.1282347.
- [3] W. Liu, H. Wang, Y. Zhang, W. Wang, and L. Qin, "I-LSH: I/O efficient c-approximate nearest neighbor search in high-dimensional space," in *Proc. IEEE 35th Int. Conf. Data Eng. (ICDE)*, Apr. 2019, pp. 1670–1673, doi: 10.1109/ICDE.2019.00169.
- [4] T. Liu, C. Rosenberg, and H. Rowley, "Clustering billions of images with large scale nearest neighbor search," in *Proc. IEEE Workshop Appl. Comput. Vis. (WACV)*, Feb. 2007, p. 28, doi: 10.1109/WACV.2007.18.
- [5] J. L. Bentley, "Multidimensional binary search trees used for associative searching," *Commun. ACM*, vol. 18, no. 9, pp. 509–517, Sep. 1975, doi: 10.1145/361002.361007.

- [6] K. L. Cheung and A. W.-C. Fu, "Enhanced nearest neighbour search on the R-tree," *ACM SIGMOD Rec.*, vol. 27, no. 3, pp. 16–21, Sep. 1998, doi: 10.1145/290593.290596.
- [7] R. Weber, H. J. Schek, and S. Blott, "A quantitative analysis and performance study for similarity-search methods in high-dimensional spaces," in *Proc. 24rd Int. Conf. Very Large Data Bases*, 1998, pp. 194–205.
- [8] O. Jafari and P. Nagarkar, "Experimental analysis of locality sensitive hashing techniques for high-dimensional approximate nearest neighbor searches," in *Proc. Australas. Database Conf.*, 2021, pp. 62–73, doi: 10.1007/978-3-030-69377-0_6.
- [9] K. Zhao, H. Lu, and J. Mei, "Locality preserving hashing," in *Proc. AAAI Conf. Artif. Intell.*, vol. 28, no. 1, Jun. 2014, doi: 10.1609/aaai.v28i1.9133.
- [10] A. Gionis, P. Indyk, and R. Motwani, "Similarity search in high dimensions via hashing," in *Proc. 25th Int. Conf. Very Large Data Bases*, 1999, pp. 518–529, Accessed: Oct. 8, 2022.
- [11] Y. Tao, K. Yi, C. Sheng, and P. Kalnis, "Efficient and accurate nearest neighbor and closest pair search in high-dimensional space," *ACM Trans. Database Syst.*, vol.
- [12] W. Liu, H. Wang, Y. Zhang, W. Wang, and L. Qin, "I-LSH: I/O efficient c-approximate nearest neighbor search in high-dimensional space," in *Proc. IEEE 35th Int. Conf. Data Eng. (ICDE)*, Apr. 2019, pp. 1670–1673, doi: 10.1109/ICDE.2019.00169.

