

# Benchmarking Nearest Neighbor Search: Influence of Local Intrinsic Dimensionality and Result Diversity in Real-World Datasets

Martin Aumüller\*

Matteo Ceccarello\*

## Abstract

This paper reconsiders common benchmarking approaches to nearest neighbor search. It is shown that the concept of local intrinsic dimensionality (LID) allows to choose query sets of a wide range of difficulty for real-world datasets. Moreover, the effect of different LID distributions on the running time performance of implementations is empirically studied. To this end, different visualization concepts are introduced that allow to get a more fine-grained overview of the inner workings of nearest neighbor search principles. The paper closes with remarks about the diversity of datasets commonly used for nearest neighbor search benchmarking. It is shown that such real-world datasets are not diverse: results on a single dataset predict results on all other datasets well.

## 1 Introduction

Nearest neighbor (NN) search is a key primitive in many computer science applications, such as data mining, machine learning and image processing. For example, Spring and Shrivastava very recently showed in [18] how nearest neighbor search methods can yield large speed-ups when training neural network models. In this paper, we study the classical  $k$ -NN problem. Given a dataset  $S \subseteq \mathbb{R}^d$ , the task is to build an index on  $S$  to support the following type of query: For a query point  $\mathbf{x} \in \mathbb{R}^d$ , return the  $k$  closest points in  $S$  under some distance measure  $D$ .

In many practical settings, a dataset consists of points represented as high-dimensional vectors. For example, word representations generated by the `glove` algorithm [17] associate with each word in a

corpus a  $d$ -dimensional real-valued vector. Common choices for  $d$  are between 50 and 300 dimensions. Finding the true nearest neighbors in such a high-dimensional space is difficult, a phenomenon often referred to as the “curse of dimensionality”. In practice, it means that finding the true nearest neighbors, in general, cannot be solved much more efficiently than by a linear scan through the dataset (requiring time  $O(n)$  for  $n$  data points) or in space that is exponential in the dimensionality  $d$ , which is impractical for large values of  $d$ .

While we cannot avoid these general hardness results, most datasets that are used in applications are not *truly* high-dimensional. This means that the dataset can be embedded into a lower-dimensional space without too much distortion. Intuitively, the intrinsic dimensionality (ID) of the dataset is the minimum number of dimensions that allows for such a representation [6]. There exist many explicit ways of find good embeddings for a given dataset. For example, the Johnson-Lindenstrauss transformation [11] allows us to embed  $n$  data points in  $\mathbb{R}^d$  into  $\Theta((\log n)/\varepsilon^2)$  dimensions such that all pairwise distances are preserved up to a  $(1+\varepsilon)$  factor with high probability. Another classical embedding often employed in practice is given by principal component analysis (PCA).

In this paper, we set our focus on the “local intrinsic dimensionality” (LID), a measure introduced by Houle in [6]. We defer a detailed discussion of this measure and its estimation to Section 2. Intuitively, the LID of a data point  $\mathbf{x}$  at a distance threshold  $r > 0$  measures how difficult it is to distinguish between points at distance  $r$  and distance  $(1 + \varepsilon)r$  in a dataset. Most importantly for this study, LID is a *local* measure that can be associated with a single query. It was stated in [7] that the LID

\*IT University of Copenhagen, Denmark, {maau,mcec}@itu.dk

might serve as a characterization of the difficulty of  $k$ -NN queries. One purpose of this paper is to shed light on this statement.

A focus of this paper is an empirical study of how the LID influences the performance of NN algorithms. To be precise, we will benchmark four different implementations [12] which employ different approaches to NN search. Three of them (HNSW [15], FAISS-IVF [10], Annoy [3]) stood out as most performant in the empirical study conducted by Aumüller et al. in [2]. Another one (ONNG) was proposed very recently [8] and shown to be competitive to the former approaches there. We base our experiments on [2] and describe their benchmarking approach and the changes we made to their system in Section 3. We analyze the LID distribution of real-world datasets in Section 4. We will see that there is a substantial difference between the LID distributions among these datasets. We will then conduct two experiments: First, we fix a dataset and choose as query set the set of points with smallest, medium, and largest estimated LIDs. As we will see, the larger the LID, the more difficult the query for all implementations. Next, we will study how the different LID distributions between datasets influence the running time performance. In a nutshell, it cannot be concluded that LID by itself is a good indicator for the relative performance of a fixed implementation over datasets. These statements will be made precise in the evaluation that is discussed in Section 5.

In the first part of our evaluation, we will work in the “classical evaluation setting of nearest neighbor search”. This means that we relate a performance measure (such as the achieved throughput measured in queries per second) to a quality measure (such as the average fraction of true nearest neighbors found over all queries). While this is the most commonly employed evaluation method, we will reason that this way of representing results in fact hides interesting details about the inner workings of an implementation. Using non-traditional visualization techniques will provide new insights into their query behavior on real-world datasets. As one example, we will see that reporting average recall on the graph-based approaches from [15, 8] hides an important detail:

For a given query, they either find all true nearest neighbors or not a single one. This behavior is not shared by the two other approaches that we consider; both yield a continuous transition from “finding no nearest neighbors” to “finding all nearest neighbors”.

Finally, we will challenge two more common benchmarking approaches to NN search. First, we will empirically observe how using small query sets (less than 1% of the dataset) affects the robustness of measurements. Second, we will discuss how *diverse* results on the considered real-world datasets are. Ideally, we want to benchmark on a collection of “interesting” datasets that show the strengths and weaknesses of individual approaches. We will conclude that there is little diversity among the considered datasets: While the individual performance observations change from dataset to dataset, the relative performance between implementations stays the same.

**1.1 Our Contributions** The main contributions of this paper are

- a detailed evaluation of the LID distribution of many real-world datasets used in benchmarking frameworks,
- an evaluation of the influence of the LID to the performance of NN search implementations,
- considerations about the diversity of results over many different datasets, and
- an exploration of different visualization techniques that shed light on individual properties of certain implementation principles.

**1.2 Related Work on Benchmarking Frameworks for NN** We use the benchmarking system described in [2] as the starting point for our study. Different approaches to benchmarking nearest neighbor search are described in [4, 5, 14]. We refer to [2] for a detailed comparison between the frameworks.

## 2 Local Intrinsic Dimensionality

We consider a distance-space  $(\mathbb{R}^d, D)$  with a distance function  $D: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ . As described in [1], we consider the distribution of distances within this

space with respect to a reference point  $\mathbf{x}$ . Such a distribution is induced by sampling  $n$  points from the space  $\mathbb{R}^d$  under a certain probability distribution. We let  $F: \mathbb{R} \rightarrow [0, 1]$  be the cumulative distribution function of distances to the reference point  $\mathbf{x}$ .

**Definition 2.1** ([6]). The local continuous intrinsic dimension of  $F$  at distance  $r$  is given by

$$\text{ID}_F(r) = \lim_{\varepsilon \rightarrow 0} \frac{\ln(F((1 + \varepsilon)r)/F(r))}{\ln((1 + \varepsilon)r/r)},$$

whenever this limit exists.

The measure relates the increase in distance to the increase in probability mass (the fraction of points that are within the ball of radius  $r$  and  $(1 + \varepsilon)r$  around the query point). Intuitively, the larger the LID, the more difficult it is to distinguish true nearest neighbors at distance  $r$  from the rest of the dataset. As described in [7], in the context of  $k$ -NN search we set  $r$  as the distance of the  $k$ -th nearest neighbor to the reference point  $\mathbf{x}$ .

**Estimating LID** We use the Maximum-Likelihood estimator (MLE) described in [13, 1] to estimate the LID of  $\mathbf{x}$  at distance  $r$ . Let  $r_1 \leq \dots \leq r_k$  be the sequence of distances of the  $k$ -NN of  $\mathbf{x}$ . The MLE  $\hat{\text{ID}}_{\mathbf{x}}$  is then

$$(2.1) \quad \hat{\text{ID}}_{\mathbf{x}} = - \left( \frac{1}{k} \sum_{i=1}^k \ln \frac{r_i}{r_k} \right)^{-1}.$$

Amsaleg et al. showed in [1] that MLE estimates the LID well.

### 3 Overview over the Benchmarking Framework

We use the `ann-benchmarks` system described in [2] to conduct our experimental study. `ann-benchmarks` is a framework for benchmarking NN search algorithms. It covers dataset creation, performing the actual experiment, and storing the results of these experiments in a transparent and easy-to-share way. Moreover, results can be explored through various plotting functionalities, e.g., by creating a website containing interactive plots for all experimental runs.

`ann-benchmarks` interfaces with a NN search implementation by calling its preprocess (index

building) and search (query) methods with certain parameter choices. Implementations are tested on a large set of parameters usually provided by the original authors of an implementation. The answers to queries are recorded as the indices of the points returned. `ann-benchmarks` stores these parameters together with further statistics such as individual query times, index size, and auxiliary information provided by the implementation. We refer to [2] for more details.

Compared to the system described in [2], we added tools to estimate the LID based on (2.1), pick “challenging query sets” according to the LID of individual points, and added new datasets and implementations. Moreover, we implemented a mechanism that allows an implementation to provide further query statistics after answering a query. To showcase this feature, all implementations considered in this study report the number of distance computations performed to answer a query.<sup>1</sup>

## 4 Algorithms and Datasets

**4.1 Algorithms** Nearest neighbor search algorithms for high dimensions are usually graph-, tree-, or hashing-based. We refer the reader to [2] for an overview over these principles and available implementations. In this study, we concentrate on the three implementations considered most performant in [2], namely HNSW [15], Annoy [3] and FAISS-IVF [10] (IVF from now on). We consider the very recent graph-based approach ONNG [8] in this study as well.

HNSW and ONNG are graph-based approaches. This means that they build a  $k$ -NN graph during the preprocessing step. In this graph, each vertex is a data point and a directed edge  $(u, v)$  means that the data point associated with  $v$  is “close” to the data point associated with  $u$  in the dataset. At query time, the graph is traversed to generate candidate points. Algorithms differ in details of the graph construction, how they build a navigation structure on top of the graph, and how the graph is traversed.

`Annoy` is an implementation of a random projec-

<sup>1</sup>We thank the authors of the implementations for their help and responsiveness in adding this feature to their library.

Dataset	Data Points	Dimensions	avg(LID)	median(LID)	Metric
SIFT [9]	1 000 000	128	21.9	19.2	Euclidean
MNIST	65 000	784	14.0	13.2	Euclidean
Fashion-MNIST [19]	65 000	784	15.6	13.9	Euclidean
GLOVE [17]	1 183 514	100	18.0	17.8	Angular/Cosine
GLOVE-2M [17]	2 196 018	300	26.1	23.4	Angular/Cosine
GNEWS [16]	3 000 000	300	21.1	20.1	Angular/Cosine

Table 1: Datasets under consideration with their average local intrinsic dimensionality (LID) computed by MLE [1] from the 100-NN of all the data points. For each dataset we use 10 000 query points.

tion forest, which is a collection of random projection trees. Each node in a tree is associated with a set of data points. It splits these points into two subsets according to a chosen hyperplane. If the dataset in a node is small enough, it is stored directly and the node is a leaf. **Annoy** employs a data-dependent splitting mechanism in which a splitting hyperplane is chosen as the one splitting two “average points” by repeatedly sampling dataset points. In the query phase, trees are traversed using a priority queue until a predefined number of points is found.

IVF builds an inverted file based on clustering the dataset around a predefined number of centroids. It splits the dataset based on these centroids by associating each point with its closest centroid. During query it finds the closest centroids and checks points in the dataset associated with those.

We remark that both IVF and HNSW implementations are the ones provided with **FAISS**<sup>2</sup>.

**4.2 Datasets** Table 1 presents an overview over the datasets that we consider in this study. We restrict our attention to datasets that are usually used in connection with Euclidean distance and Angular/Cosine distance. For each dataset, we compute the LID distribution with respect to the 100-NN as discussed in Section 2. SIFT, MNIST, and GLOVE are among the most-widely used datasets for benchmarking nearest neighbor search algorithms. Fashion-MNIST is considered as a replacement for MNIST, which is usually considered too easy for machine learning tasks [19].

Figure 1 provides a visual representation of the estimated LID distribution of each dataset, for

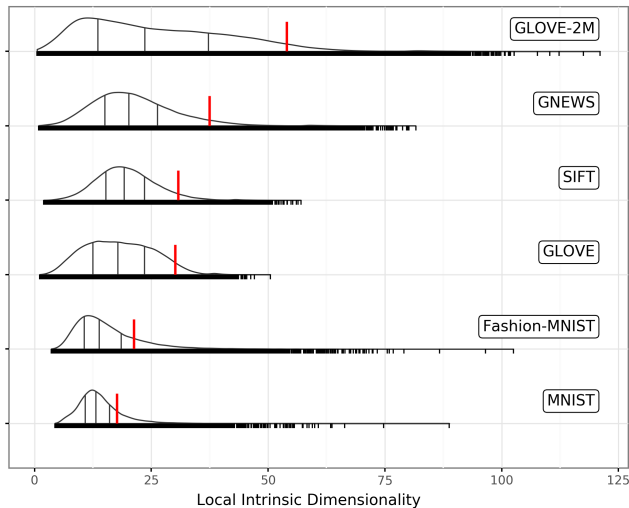


Figure 1: LID distribution for each dataset. Ticks below the distribution curves represent single queries. Lines within each distribution curve correspond to the 25, 50 and 75 percentile. The red line marks the 10 000 largest estimated LID, which we use as a threshold value to define *hard* query sets.

$k = 1000$ . While the datasets differ widely in their original dimensionality, the median LID ranges from around 13 for MNIST to about 23 for GLOVE-2M. The distribution of LID values is asymmetric and shows a long tail behavior. MNIST, Fashion-MNIST, SIFT, and GNEWS are much more concentrated around the median compared to the two *glove*-based datasets.

## 5 Evaluation

This section reports on the results of our experiments. Due to space constraints, we only present some selected results. More results can be explored

<sup>2</sup><https://github.com/facebookresearch/faiss>

via interactive plots at <http://ann-benchmarks.com/edml19/>, which also contains a link to the source code repository. For a fixed implementation, the plots presented here consider the pareto frontier over all parameter choices [2]. Tested parameter choices and the associated plots are available on the website.

**Experimental Setup** Experiments were run on 2x 14-core Intel Xeon E5-2690v4 (2.60GHz) with 512GB of RAM using Ubuntu 16.10 with kernel 4.4.0. Index building was multi-threaded, queries were answered sequentially in a single thread.

**Quality and Performance Metrics** As quality metric we measure the individual recall of each query, i.e., the fraction of points reported by the implementation that are among the true  $k$ -NN. As performance metric, we record individual query times and the total number of distance computations needed to answer all queries. We usually report on the throughput (the average number of queries that can be answered in one second, in the plots denoted as QPS for *queries per second*), but we will also inspect individual query times.

**Objectives of the Experiments** Our experiments are tailored to answer the following questions:

- (Q1) How robust are measurements when splitting query set and dataset at random multiple times? (Section 5.1)
- (Q2) How does the LID of a query set influence the running time performance? (Section 5.2)
- (Q3) How diverse are measurements obtained on datasets? Do relative differences between the performance of different implementations stay the same over multiple datasets? (Section 5.3)
- (Q4) How well does the number of distance computations reflect the relative running time performance of the tested implementations? (Section 5.3)
- (Q5) How concentrated are quality and performance measures around their mean for the tested implementations? (Section 5.4)

**Choosing Query Sets** All runs are based on 10 000 queries chosen from the dataset. Depending on the question we try to answer, we employ different selection strategies for choosing the query set.

To answer Q1 we employ random splits as done in the cross-validation technique: 10 000 randomly selected points are used as queries, while the rest are used as data points. Repeating the procedure 10 times yields 10 queryset/dataset pairs.

For the discussion of Q2–Q5 we select three different querysets from each dataset: (i) The dataset that contains the 10 000 points with the lowest estimated LID (which we denote *easy*), 10 000 points around the data point with median estimated LID (denoted *medium*), and the 10 000 points with the largest estimated LID (dubbed *hard*). Figure 1 marks with a red line the LID used as a threshold to build the *hard* queryset.

**5.1 Robustness of Random Split** Figure 2 shows the result of ten cross-validation runs on random splits of GLOVE for ONNG and Annoy. The plot relates the average number of queries that can be answered in one second to the average recall achieved with a certain parameter configuration. This means that we want data points to be up and to the right. The plot shows that the 10 splits give very similar results and represent the performance of the implementation accurately. It is interesting to see that a query set representing less than 1% of the dataset points shows such a robust behavior. We conclude that the random split allows for robust measurements.

**5.2 Influence of LID on Performance** Figure 3 shows results for the influence of using only points with low, middle, and large estimated LID as query points, in SIFT and GLOVE-2M. We observe a clear influence of the LID of the query set on the performance: the larger the LID, the more down and to the left the graphs move. This means that, for higher LID, it is more expensive, in terms of time, to answer queries with good recall. For all datasets except GLOVE-2M, all implementations were still able to achieve close to perfect recall with the parameters set. This means that all but one

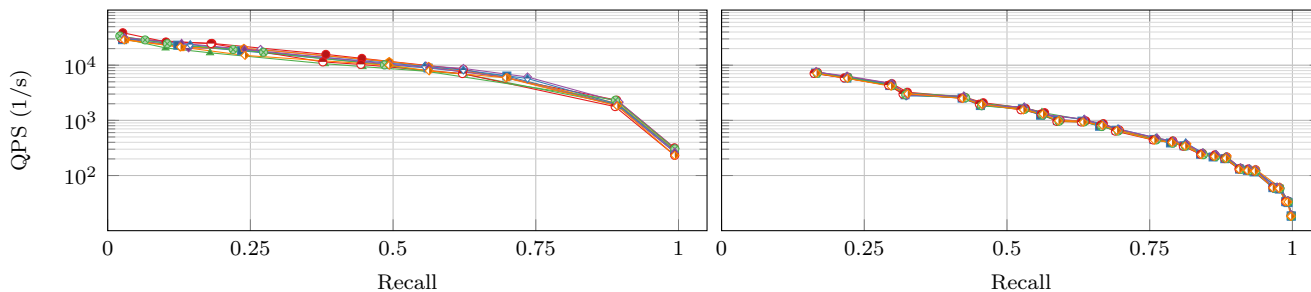


Figure 2: (Approximate) Recall-QPS (1/s) tradeoff - up and to the right is better. 10 random splits of the dataset; left: ONNG, right: ANNOY.

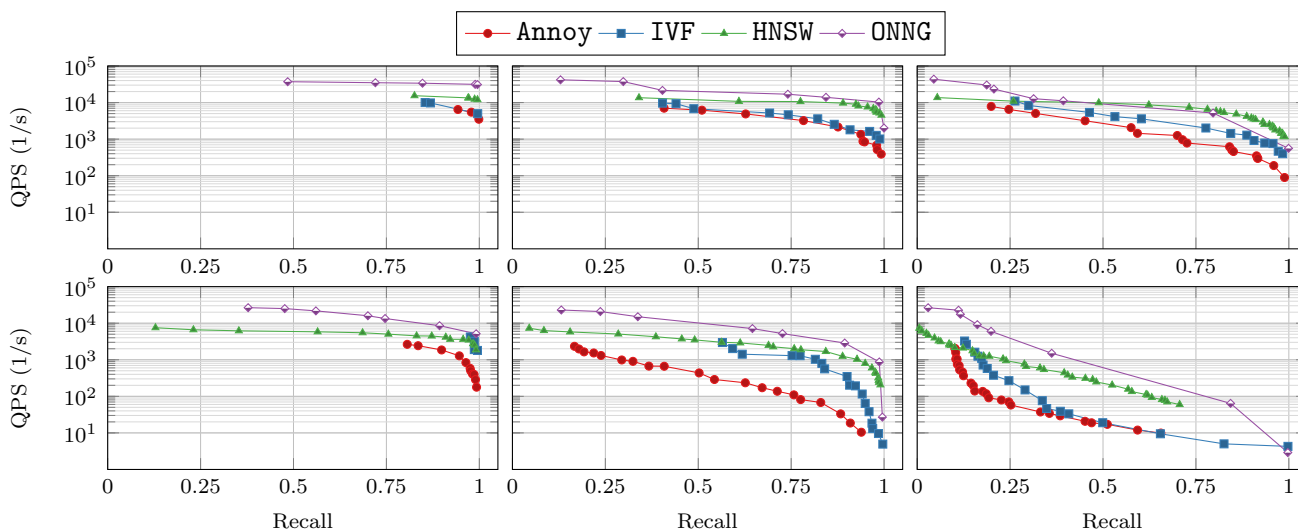


Figure 3: (Approximate) Recall-QPS (1/s) tradeoff - up and to the right is better; top: SIFT, bottom: GLOVE-2M; left: easy, middle: middle, right: hard.

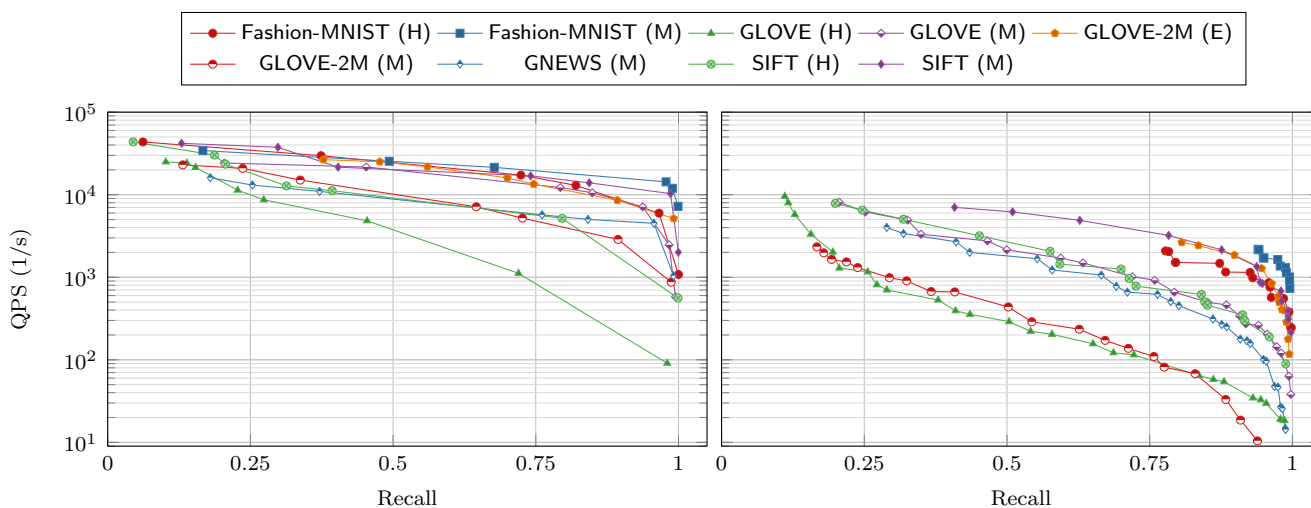


Figure 4: (Approximate) Recall-QPS (1/s) tradeoff - up and to the right is better; left: ONNG, right: Annoy; (E) — easy, (M) — medium, (H) — hard.

of the tested datasets do not contain too many “noisy queries”. Already the queries around the median prove challenging for most implementations. For the most difficult queries (according to LID), only IVF and ONNG achieve close to perfect recall on GLOVE-2M.

Figure 4 reports on the results of ONNG and Annoy on selected datasets. Comparing results to the LID measurements depicted in Figure 1, the estimated median LID gives a good estimate on the relative performance of the algorithms on the datasets. As an exception, SIFT (M) is much easier than predicted by its LID distribution. In particular for Annoy, the hard SIFT instance is as challenging as the medium GLOVE version. The easy version of GLOVE-2M turns out to be efficiently solvable by both implementations (taking about the same time as it takes to answer the hard instance of Fashion-MNIST, which has a much higher LID). From this, we cannot conclude that LID as a single indicator explains performance differences of an implementation on different datasets. However, more careful experimentation is needed before drawing a final conclusion. In our setting, the LID estimation is conducted for  $k = 1000$ , while queries are only searching for the 10 nearest neighbors. Moreover, the estimation using MLE might not be accurate enough on these datasets. We leave the investigation of these two directions as future work.

As a side note, we remark that Fashion-MNIST is as difficult to solve as MNIST for all implementations, and is by far the easiest dataset for all implementations. Thus, while there is a big difference in the difficulty of solving the classification task [19], there is no measurable difference between these two datasets in the context of NN search.

**5.3 Diversity of Results** Figure 5 gives an overview over how algorithms compare to each other among all “medium difficulty” datasets. We consider two metrics, namely the number of queries per second (top plot), and the number of distance computations (bottom plot). For two different average recall thresholds (0.75 and 0.9) we then select, for each algorithm, the best performing parameter configuration that attains at least that recall. For each dataset, the plots report the ratio

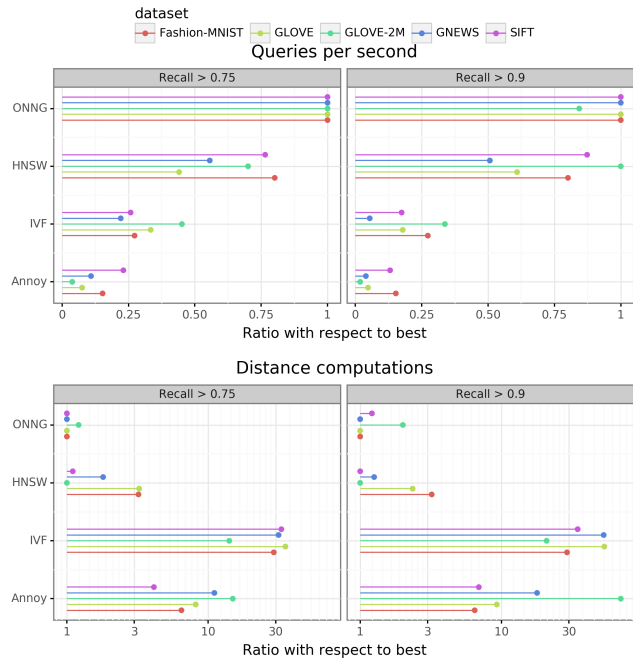


Figure 5: Ranking of algorithms on five different datasets, according to recall  $\geq 0.75$  and  $\geq 0.9$ , and according to two different performance measures: number of queries per second (top) and number of distance computations (bottom). Both plots report the ratio with the best performing algorithm on each dataset: for the queries per second metric a larger ratio is better, for the number of distance computations metric, a smaller ratio is better.

with the best performing algorithm on that dataset, therefore the best performer is reported with ratio 1. Considering different dataset, we see that there is little variation in the ranking of the algorithms. Only the two graph-based approaches trade ranks, all other rankings are stable. Interestingly, Annoy makes much fewer distance computations but is consistently outperformed by IVF.<sup>3</sup>

Comparing the number of distance computations to running time performance, we see that an increase in the number of distance computations is not reflected in a proportional decrease in the number of queries per second. This means that the

<sup>3</sup>We note that IVF counts the initial comparisons to find the closest centroids as distance computations, whereas Annoy did not count the inner product computations during tree traversal. We plan to report on updated results in the workshop.



candidate set generation is in general more expensive for graph-based approaches, but the resulting candidate set is of much higher quality and fewer distance computations have to be carried out. Generally, both graph-based algorithms are within a factor 2 from each other, whereas the other two need much larger candidate lists to achieve a certain recall. The relative difference usually ranges from 5x to 30x more distance computations for the non-graph based approaches, in particular at high recall. This translates well into the performance differences we see in this setting: consider for instance Figure 3, where the lines corresponding to HNSW and ONNG upper bound the lines relative to the other two algorithms.

#### 5.4 Reporting the distribution of performance

In the previous sections, we made extensive use of recall/queries per second plots, where each configuration of each algorithm results in a single point, namely the average recall and the inverse of the average query time. As we shall see in this section, concentrating on averages can hide interesting information in the context of  $k$ -NN queries. In fact, not all queries are equally difficult to answer. Consider the plots in Figure 6, which report performance of the four algorithms<sup>4</sup> on the GLOVE-2M dataset, medium difficulty. The left plot reports the recall versus the number of queries per second, and black dots correspond to the averages. Additionally, for each configuration, we report the distribution of the recall scores: the baseline of each recall curve is positioned at the corresponding queries per second performance. Similarly, the right plot reports on the inverse of the individual query times (the average of these is the QPS in the left plot) against the average recall. In both plots, the best performance is achieved towards the top-right corner.

Plotting the distributions, instead of just reporting the averages, uncovers some interesting behaviour that might otherwise go unnoticed, in particular with respect to the recall. The average recall gradually shifts towards the right as the effect of more and more queries achieving good recalls. Per-

haps surprisingly, for graph-based algorithms this shift is very sudden: most queries go from having recall 0 to having recall 1, taking no intermediate values. Taking the average recall as a performance metric is convenient in that it is a single number to compare algorithms with. However, the same average recall can be attained with very different distributions: looking at such distributions can provide more insight.

For the plot on the right, we observe that individual query times of Annoy, IVF, and ONNG are well concentrated around their mean. However, for HNSW query times fluctuate widely around their mean; in fact, the average query time is not well reflected in the query time distribution.

For space reasons, we do not report other parameter configurations and datasets, which nonetheless show similar behaviours.

## 6 Summary

In this paper we studied the influence of LID to the performance of nearest neighbor search algorithms. We showed that LID allows to choose query sets of a wide range of difficulty from a given dataset. We also showed how different LID distributions influence the running time performance of the algorithms. In this respect, we could not conclude that the LID alone can predict running time differences well. In particular, SIFT is usually easier for the algorithms, while GLOVE’s LID distribution would predict it to be the easier dataset of the two.

We introduced novel visualization techniques to show the uncertainty within the answer to a set of queries, which made it possible to show a clear difference between the graph-based algorithms and the other approaches.

We hope that this study initiates the search for more diverse datasets, or for theoretical reasoning why certain algorithmic principles are generally better suited for nearest neighbor search. From a more practical side, we would be interested in seeing whether the LID can be used in the design of NN algorithms to guide the search process or parameter selection. For example, algorithms could adapt to the LID of the current set of  $k$ -NN found so far, stopping early or spending more time depending on the LID of the candidates.

<sup>4</sup>In order not to clutter the plots, we fixed parameters as follows: IVF — number of lists 8192; Annoy — number of trees 100; HNSW — efConstruction 500; ONNG — edge 100, outdegree 10, indegree 120.



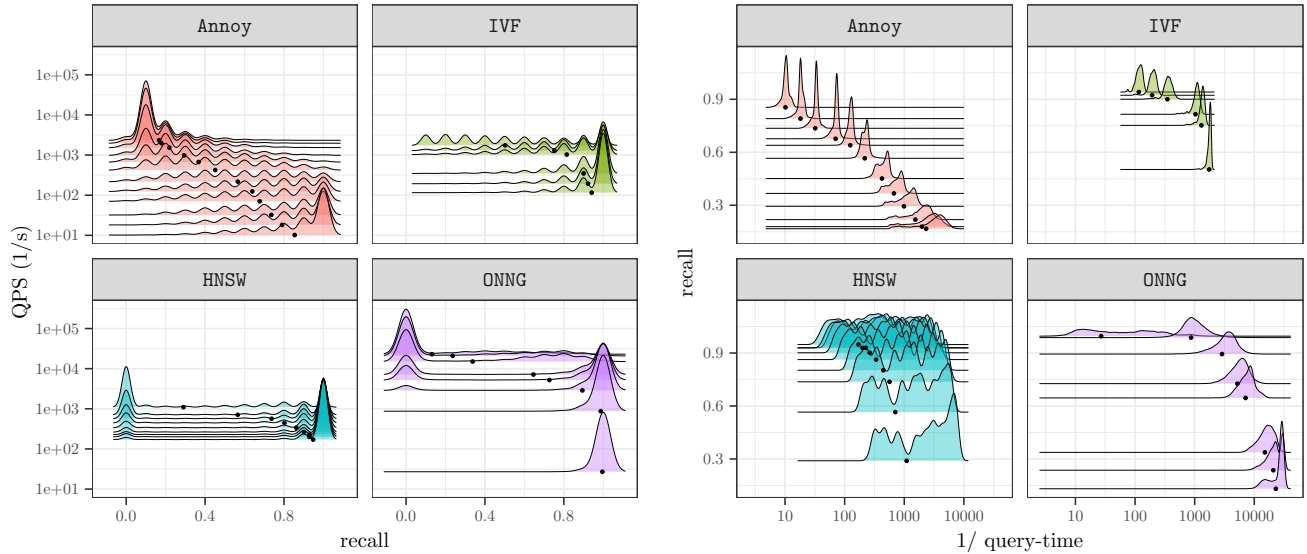


Figure 6: Distribution of performance for queries on the GLOVE-2M (medium difficulty) dataset. Looking just at the average performance can hide interesting behaviour.

## References

- [1] Amsaleg, L., Chelly, O., Furon, T., Girard, S., Houle, M.E., Kawarabayashi, K.I., Nett, M.: Estimating local intrinsic dimensionality. In: KDD'15. pp. 29–38. ACM (2015)
- [2] Aumüller, M., Bernhardsson, E., Faithfull, A.: ANN-Benchmarks: A Benchmarking Tool for Approximate Nearest Neighbor Algorithms. In: SISAP'17. pp. 34–49 (2017)
- [3] Bernhardsson, E.: Annoy, <https://github.com/spotify/annoy>
- [4] Curtin, R.R., Cline, J.R., Slagle, N.P., March, W.B., Ram, P., Mehta, N.A., Gray, A.G.: ML-PACK: A scalable C++ machine learning library. *J. of Machine Learning Research* 14, 801–805 (2013)
- [5] Edel, M., Soni, A., Curtin, R.R.: An automatic benchmarking system. In: NIPS 2014 Workshop on Software Engineering for Machine Learning (2014)
- [6] Houle, M.E.: Dimensionality, discriminability, density and distance distributions. In: Data Mining Workshops (ICDMW). pp. 468–473. IEEE (2013)
- [7] Houle, M.E., Schubert, E., Zimek, A.: On the correlation between local intrinsic dimensionality and outlieriness. In: SISAP'18. pp. 177–191 (2018)
- [8] Iwasaki, M., Miyazaki, D.: Optimization of Indexing Based on k-Nearest Neighbor Graph for Proximity Search in High-dimensional Data. ArXiv e-prints (Oct 2018)
- [9] Jégou, H., Douze, M., Schmid, C.: Product quantization for nearest neighbor search. *IEEE Trans. Pattern Anal. Mach. Intell.* 33(1), 117–128 (2011)
- [10] Johnson, J., Douze, M., Jégou, H.: Billion-scale similarity search with gpu. *CoRR abs/1702.08734* (2017)
- [11] Johnson, W.B., Lindenstrauss, J., Schechtman, G.: Extensions of lipschitz maps into banach spaces. *Israel Journal of Mathematics* 54(2), 129–138 (1986)
- [12] Kriegel, H., Schubert, E., Zimek, A.: The (black) art of runtime evaluation: Are we comparing algorithms or implementations? *Knowl. Inf. Syst.* 52(2), 341–378 (2017)
- [13] Levina, E., Bickel, P.J.: Maximum likelihood estimation of intrinsic dimension. In: NIPS'15. pp. 777–784 (2005)
- [14] Li, W., Zhang, Y., Sun, Y., Wang, W., Zhang, W., Lin, X.: Approximate nearest neighbor search on high dimensional data - experiments, analyses, and improvement (v1.0). *CoRR abs/1610.02455* (2016)
- [15] Malkov, Y.A., Yashunin, D.A.: Efficient and robust approximate nearest neighbor search using Hierarchical Navigable Small World graphs. ArXiv e-prints (Mar 2016)
- [16] Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. *CoRR abs/1301.3781* (2013)
- [17] Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: Empirical Methods in Natural Language Processing (EMNLP). pp. 1532–1543 (2014)

- [18] Spring, R., Shrivastava, A.: Scalable and sustainable deep learning via randomized hashing. In: KDD'17. pp. 445–454 (2017)
- [19] Xiao, H., Rasul, K., Vollgraf, R.: Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. CoRR abs/1708.07747 (2017)