

VIBE: Vector Index Benchmark for Embeddings

Elias Jääsaari* Ville Hyvönen† Matteo Ceccarello§ Teemu Roos* Martin Aumüller‡
 *University of Helsinki †Aalto University §University of Padova ‡IT University of Copenhagen

 [Code: https://github.com/vector-index-bench/vibe](https://github.com/vector-index-bench/vibe)

 [Website: https://vector-index-bench.github.io](https://vector-index-bench.github.io)

Abstract

Approximate nearest neighbor (ANN) search is a performance-critical component of many machine learning pipelines. Rigorous benchmarking is essential for evaluating the performance of vector indexes for ANN search. However, the datasets of the existing benchmarks are no longer representative of the current applications of ANN search. Hence, there is an urgent need for an up-to-date set of benchmarks. To this end, we introduce Vector Index Benchmark for Embeddings (VIBE), an open source project for benchmarking ANN algorithms. VIBE contains a pipeline for creating benchmark datasets using dense embedding models characteristic of modern applications, such as retrieval-augmented generation (RAG). To replicate real-world workloads, we also include out-of-distribution (OOD) datasets where the queries and the corpus are drawn from different distributions. We use VIBE to conduct a comprehensive evaluation of SOTA vector indexes, benchmarking 21 implementations on 12 in-distribution and 6 out-of-distribution datasets.

1 Introduction

In modern machine learning applications, data, such as text and images, are often stored as *embeddings* that are generated by deep learning models. Embeddings are continuous vector representations that have the property that semantically similar items have similar representations in the embedding space. This makes efficient similarity search in high-dimensional vector spaces an essential operation in applications where corpus items that are semantically similar to the query item are retrieved, such as recommendation systems (Covington et al., 2016; Yang et al., 2020), information retrieval (Xiong et al., 2021; Long et al., 2022), and question answering (Seo et al., 2019; Lewis et al., 2021).

Since datasets in these applications are typically large and high-dimensional, *approximate nearest neighbor* (ANN) search, where *vector indexes* are used to perform approximate k -nn queries, is often used to ensure fast response times. An important recent use case is retrieval-augmented generation (RAG) (Guu et al., 2020; Lewis et al., 2020; Borgeaud et al., 2022; Li et al., 2024; Wang et al., 2024).

Real-world workloads for ANN search are often *out-of-distribution* (OOD): the corpus and the queries have significantly different distributions. For instance, in multi-modal search, the corpus may consist of images and the queries may be given as text. Another recent OOD use case is approximate attention computation in LLMs (Bertsch et al., 2023; Zandieh et al., 2023; Liu et al., 2024; Mazare et al., 2025), where the task is to find the keys that have the largest inner products with the query.

Because of the significance of ANN search in modern AI applications, accurate performance evaluation of vector indexes is essential. Currently, the ANN-benchmarks project (Aumüller et al., 2020) includes the most comprehensive collection of ANN algorithms. However, most of its datasets, such as MNIST, Fashion-MNIST, and SIFT, are not representative of modern applications since they have representations, such as raw pixels and image descriptors, that have been superseded by dense vector embeddings. In addition, there does not exist a systematic benchmark for the recent important OOD

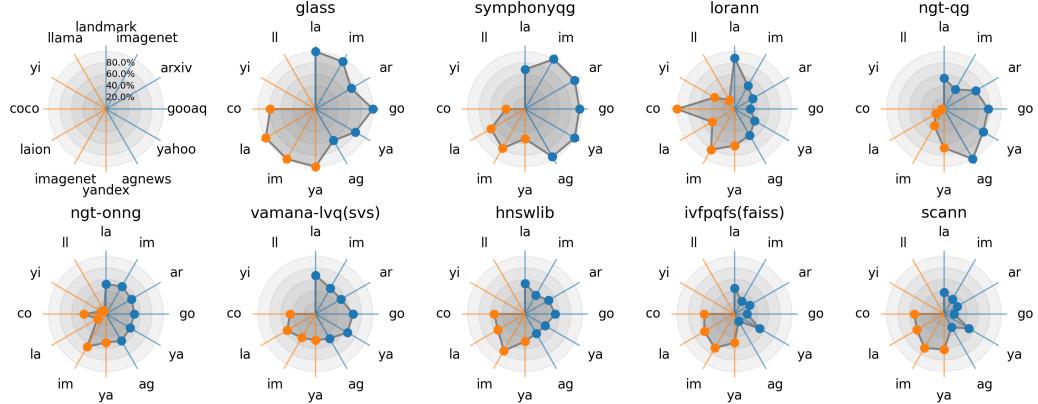


Figure 1: Algorithm performance of the fastest configuration with average recall $\geq 95\%$ on datasets with in-distribution and out-of-distribution queries, in terms of queries per second relative to the best algorithm on each dataset. Each circle corresponds to an algorithm, arranged by decreasing average rank on the datasets. Datasets are arranged in clockwise order by decreasing difficulty (as measured by the *relative contrast* (He et al., 2012)), with `llama` being the easiest. Axes with a missing dot correspond to datasets where the algorithm could not reach 95% recall. A missing axis for the `llama` and `yi` datasets indicates that the corresponding algorithm does not support inner product similarity.

use cases of ANN search, such as multi-modal search and approximate attention computation. Hence, there is an urgent need for a comprehensive and up-to-date set of benchmarks whose datasets consist of (i) modern embedding datasets and (ii) OOD datasets.

To address these issues, we introduce Vector Index Benchmark for Embeddings (VIBE), a benchmarking suite for ANN algorithms on modern embedding datasets. VIBE contains a straightforward pipeline for creating benchmark datasets: we use popular embedding models to generate representative embedding datasets from sources such as ArXiv and ImageNet. We also introduce novel approximate attention computation benchmark datasets. The benchmark is open, ongoing, and extensible, enabling authors of new ANN algorithms to easily integrate their implementations. Thus, beyond presenting an up-to-date performance evaluation, our work enables rigorous evaluation in the future. VIBE features an interactive website that enables deeper analysis of results, for instance assessing the robustness of methods by viewing the performance breakdown for easy and hard queries.

Key takeaways. We use VIBE to evaluate the performance of state-of-the-art ANN algorithms. An overview of our evaluation for selected algorithms is presented in Figure 1. Our key findings are: (i) The top-performing graph-based (SymphonyQG, Glass, NGT-QG) and clustering-based (LoRANN, ScaNN, IVF-PQ) algorithms utilize different quantization methods. (ii) On text-to-image OOD datasets, the best regular methods, although their performance deteriorates significantly on OOD queries (see Figure 6), still, somewhat surprisingly, outperform specialized OOD methods (see Figure 7). (iii) On approximate attention computation datasets, regular ANN methods fail to reach the highest recall levels, and specialized OOD methods outperform them (see Figure 8).

We summarize our main contributions below:

- **Benchmark for modern ML use cases.** VIBE is a benchmarking framework for evaluating the performance of vector indexes on embedding datasets used in modern applications, such as RAG.
- **Extensible embedding pipeline.** VIBE contains a pipeline for generating new benchmark datasets using recent embedding models, enabling effortless updating of the datasets.
- **Comprehensive evaluation.** We use VIBE to conduct a comprehensive performance evaluation of SOTA vector indexes, comparing 21 implementations on 12 in-distribution datasets.
- **Evaluation in OOD setting.** We use VIBE to conduct a systematic performance evaluation of both regular ANN methods and specialized OOD methods for out-of-distribution queries on 6 datasets.
- **Fine-grained performance metrics.** In addition to the average recall, VIBE measures more fine-grained performance metrics, such as a performance breakdown for easy/hard queries and in-distribution/out-of-distribution queries.

2 Background

2.1 ANN search

Definition. Denote by $\{\mathbf{c}_1, \dots, \mathbf{c}_m\} \subset \mathbb{R}^d$ the *corpus* and let $d : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be a dissimilarity measure. The task of k -nearest neighbor (k -nn) search is to retrieve the k corpus points that are most similar to the *query point* $\mathbf{x} \in \mathbb{R}^d$, i.e., to retrieve

$$\text{NN}_k(\mathbf{x}) = \{j \in [m] : d(\mathbf{x}, \mathbf{c}_j) \leq d(\mathbf{x}, \mathbf{c}_{(k)})\},$$

where we denote by $\mathbf{c}_{(1)}, \dots, \mathbf{c}_{(m)}$ the corpus points ordered w.r.t. their dissimilarity to \mathbf{x} in ascending order. The task of *approximate nearest neighbor* (ANN) search is to design a *vector index* that enables approximating the exact k -nn solution in sublinear time. In this paper, we use the term *approximate* to refer to *inexact solutions* without guarantees (Aumüller and Ceccarello, 2023).

Dissimilarity measures. The most common dissimilarity measures for ANN search are the Euclidean distance $d(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|_2$, the cosine distance $d(\mathbf{a}, \mathbf{b}) = 1 - \langle \mathbf{a}/\|\mathbf{a}\|_2, \mathbf{b}/\|\mathbf{b}\|_2 \rangle$, and the negative inner product $d(\mathbf{a}, \mathbf{b}) = -\langle \mathbf{a}, \mathbf{b} \rangle$. The special case where $d(\mathbf{a}, \mathbf{b}) = -\langle \mathbf{a}, \mathbf{b} \rangle$ is often called *maximum inner product search* (MIPS). Embeddings are commonly normalized to unit vectors, in which case all three dissimilarity measures yield the same nearest neighbors.

Performance measures. The effectiveness of ANN algorithms is typically measured by the average *recall*, i.e., the fraction of returned neighbors for a given query point that are among the true k nearest neighbors of that point. The efficiency of ANN algorithms is typically measured by the average query time, or, equivalently, by throughput, which is measured by *queries per second* (QPS).

Difficulty measures. To assess the difficulty of a query \mathbf{x} we consider its *relative contrast* at k , $RC_k(\mathbf{x}) = \frac{\sum_{j \in [m]} d(\mathbf{x}, \mathbf{c}_j)/m}{d(\mathbf{x}, \mathbf{c}_{(k)})}$, that is the ratio of the average distance from the query to the dataset with the distance to its k -th nearest neighbor (He et al., 2012). The relative contrast is a good proxy for the effort required to distinguish close and far away points (Aumüller and Ceccarello, 2021), with small values characterizing difficult queries.

Out-of-distribution (OOD) setting. There is a recent interest for ANN search in the *out-of-distribution* (OOD) setting (e.g., Cayton and Dasgupta, 2007; Jaiswal et al., 2022; Hyvönen et al., 2022; Chen et al., 2024; Jääsaari et al., 2024; Mazaré et al., 2025), where the distributions of the corpus points and the query points are different. If there is a large difference between the query distribution and the corpus distribution, using a sample from the actual query distribution to adapt the index structure can speed up ANN search significantly (Jaiswal et al., 2022; Hyvönen et al., 2022; Chen et al., 2024). The OOD setting is common in multi-modal search, where the corpus and the queries have different modalities—for instance, the queries are given as text, and the corpus consists of images—even though they are embedded in the same vector space. Another recent example is approximate attention computation in long-context LLMs (Liu et al., 2024; Mazaré et al., 2025).

2.2 Taxonomy of ANN algorithms

Graphs. Graph-based methods, such as HNSW (Malkov and Yashunin, 2018) and SymphonyQG (Gou et al., 2025) use an approximate k -nn graph of the corpus as an index structure and retrieve the approximate nearest neighbors of the query point via greedy graph traversal/beam search.

Clustering. Clustering-based methods cluster the corpus and retrieve nearest neighbors by evaluating the distances from the query point only to the corpus points in the clusters closest to the query point. Before exact distance evaluation, this set of points is typically further pruned by a scoring method, such as product quantization (PQ) in IVF-PQ (Jegou et al., 2011), anisotropic quantization in ScaNN (Guo et al., 2020), or reduced-rank regression in LoRANN (Jääsaari et al., 2024).

Trees. Tree-based methods utilize space-partitioning trees, such as k -d trees (Friedman et al., 1977), principal component trees (Sproull, 1991), or random projection trees (Dasgupta and Freund, 2008) as index structures. The most efficient tree-based methods, such as Annoy (Spotify, 2013) and MRPT (Hyvönen et al., 2016) use ensembles of randomized trees.

Hashing. Hashing-based methods, such as PUFFINN (Aumüller et al., 2019) and FALCONN++ (Pham and Liu, 2022), partition the space by using locality-sensitive hash functions and evaluate the distances from the query point to only those corpus points that receive a similar hash value to the query. These approaches provide strong probabilistic correctness guarantees.

Table 1: VIBE in-distribution embedding datasets. For references and detailed data sources and embedding models used, see Appendix B. Embeddings with distance marked as “any” are normalized to unit norm, and each algorithm can choose an appropriate distance measure.

Data Source	Model	Type	n	d	Distance
AGNews	MXBAI	Text	769 382	1024	euclidean
arXiv abstracts	Nomic Text	Text	1 344 643	768	any
Google Q&A	DistilRoBERTa	Text	1 475 024	768	any
ImageNet	CLIP	Image	1 281 167	512	any
Landmark	Nomic Vision	Image	760 757	768	any
Yahoo Answers	MiniLM	Text	677 305	384	any
CelebA	ResNet	Image	201 599	2048	cosine
CCNews	Nomic Text	Text	495 328	768	any
CodeSearchNet	Jina	Code	1 374 067	768	cosine
-	GloVe	Word	1 192 514	200	cosine
Landmark	DINOv2	Image	760 757	768	cosine
Simple Wiki	OpenAI	Text	260 372	3072	any

3 Related work

Benchmarks. The ANN-benchmarks project¹ (Aumüller et al., 2020) is the most well-known and comprehensive benchmarking suite for vector indexes. However, its datasets are no longer representative of the datasets used in the modern applications of ANN search, and the project does not contain a systematic evaluation of the OOD setting. The big-ann-benchmarks project (Simhadri et al., 2022, 2024) has single-dataset tracks for benchmarking ANN algorithms on specialized tasks such as filtered search, streaming search, and OOD search. Some vector database companies have their own benchmarks, such as VectorDBBench², but they benchmark industrial vector database products, in contrast to open source ANN search implementations.

Performance evaluations. Earlier performance evaluations (Li et al., 2019; Aumüller et al., 2020) do not contain modern embedding datasets or the recent top-performing ANN methods. More recent performance evaluations (Wang et al., 2021; Azizi et al., 2025) consider only graph-based ANN methods and do not provide an open and extensible benchmarking framework.

4 VIBE

4.1 Datasets

In-distribution datasets. The in-distribution datasets included in VIBE are listed in Table 1. We create embedding datasets by using widely-used text and image embedding models to embed common text and image data sources. Our text embedding models cover both the relatively older but widely applied models DistilRoBERTa (Sanh et al., 2019; Zhuang et al., 2021) and MiniLM (Wang et al., 2020), as well as the recent top-performing industry models nomic-embed-text-v1.5 (Nussbaum et al., 2025), mxbai-embed-large-v1 (Lee et al., 2024; Li and Li, 2023), and text-embedding-3-large (OpenAI, 2024). Similarly, for images, we use embeddings from a ResNet-50 He et al. (2016) convolutional neural network, as well as more powerful recent image embeddings CLIP (Radford et al., 2021), DINOv2 (Oquab et al., 2024), and Nomic AI’s nomic-embed-vision-v1.5 (Nussbaum et al., 2024). Finally, we also include the older GloVe (Pennington et al., 2014) word embedding dataset due to its prevalence in earlier benchmarks. These datasets exhibit a diverse range of difficulty measured by relative contrast; see Appendix C.

Each text embedding model is used with a text source that is appropriate for the model’s context length. For most text embeddings, the model outputs are normalized to unit norm by default. In this case, cosine similarity, Euclidean distance, and inner product are equivalent, and each algorithm can choose an appropriate distance measure. Otherwise, we use the Euclidean distance to represent a clustering task. For image embeddings that are not normalized, we use cosine similarity.

¹<https://github.com/erikbern/ann-benchmarks>

²<https://github.com/zilliztech/VectorDBBench>

Table 2: VIBE out-of-distribution datasets. For detailed data sources and embedding models used, see Appendix B. The *yi* and *llama* datasets use the inner product as the dissimilarity measure as appropriate for approximate attention computation.

Data Source	Model	Type	<i>n</i>	<i>d</i>	Distance
MS COCO	Nomic Text/Vision	Text-to-Image	282 360	768	any
ImageNet (+captions)	ALIGN	Text-to-Image	1 281 167	640	any
LAION	CLIP	Text-to-Image	1 000 448	512	any
Yandex	SE-ResNeXt/DSSM	Text-to-Image	1 000 000	200	cosine
Yi-6B-200K	-	Attention	187 843	128	IP
Llama-3-8B-262k	-	Attention	256 921	128	IP

Out-of-distribution datasets. To represent multimodal search applications, we include four text-to-image datasets in VIBE (see Table 2). We use the ALIGN (Jia et al., 2021) and the nomic-embed-text/vision (Nussbaum et al., 2024, 2025) models to create two new datasets, and we also include subsets of the pre-existing LAION and Yandex datasets that have been used previously to compare OOD algorithms (Simhadri et al., 2024; Chen et al., 2024). We also include two datasets for approximate attention computation in transformer models (Liu et al., 2024; Mazaré et al., 2025) by extracting keys and queries from the attention computation of two long-context large language models, Yi-6B-200K and Llama-3-8B-Instruct-262k. For details, see Appendix B.

For all OOD datasets, we include a larger training sample from the query distribution to enable the development of new algorithms for the OOD setting. Currently, only the recent MLANN (Hyvönen et al., 2022), LoRANN (Jääsaari et al., 2024), and RoarGraph (Chen et al., 2024) methods use these query samples. MLANN and RoarGraph also require as input the nearest neighbors of the training queries from the corpus which we precompute and provide with the datasets.

For the OOD datasets, we use a separate sample of 1000 points from the query distribution as the test queries. For in-distribution datasets, we sample 1000 points from the dataset and use the rest of the points as the corpus to build the vector indexes.

Quantized datasets. State-of-the-art embeddings are increasingly trained to perform well when quantized to 8-bit integer or even binary precision³ due to their high computational and storage cost. For example, distance computations using the Hamming distance on binarized data are much faster to evaluate using specialized CPU instructions. VIBE contains several embeddings that quantize well, functions to quantize them, and a subset of our default algorithms support these quantized datasets.

4.2 Support for research

VIBE features additional components that enable future research on vector indexes.

Broad hardware support. VIBE builds on the performance evaluation code developed in ANN-benchmarks (Aumüller et al., 2020), extending it with support for a broad range of hardware platforms. Dropping the dependency on Docker, our pipeline supports high-performance computing (HPC) platforms by using Apptainer for containerization. Further, the benchmark also supports GPU-based implementations (see Figure 5, right panel).

Visualization tools. Our benchmark features a companion website⁴ that provides interactive plots (see Appendix A) to further explore the results and characteristics of the data. In particular, the website allows comparing algorithms at different recall levels, and to interactively explore the performance of any configuration of each algorithm, even the ones not on the Pareto frontier.

Open and extendable. The benchmark suite and website are designed to easily accommodate extensions by adding both new datasets and algorithms. The benchmark is open source and available at <https://github.com/vector-index-bench/vibe>.

³see e.g. <https://huggingface.co/blog/embedding-quantization>

⁴<https://vector-index-bench.github.io/>

Table 3: Summary of the evaluated algorithms. ‘*’ marks implementations that have not been included in earlier benchmarks. See Table 6 in Appendix D for citations and implementation details.

Graphs			Trees	Clustering	Hashing
Glass	PyNNDescent	LVQ*	ANNOY	IVF-PQ	PUFFINN
NGT-ONNG	HNSW	LeanVec*	MRPT	IVF	FALCONN++*
NGT-QG	NSG	RoarGraph*	MLANN*	ScaNN	
Vamana	FlatNav*	SymphonyQG*		LoRANN*	

5 Performance evaluation

5.1 Experimental setup

Each experiment is run on a compute node with 384 GB of RAM and two Intel Xeon Gold 6230 (Cascade Lake) CPUs. Each CPU has 20 cores and support for AVX-512 instructions. All algorithms are benchmarked using a single core, but we build multiple vector indexes in parallel.

By default, we use $k = 100$ because it is representative of real-world applications: retrieval systems commonly search for a larger number of nearest neighbors that are then post-processed e.g. by filtering, enforcing diversity, or using more expensive reranker models. We use recall to measure the effectiveness and queries per second (QPS) to measure the efficiency of the algorithms.

Algorithms. We survey the recent literature and the earlier benchmarking efforts (e.g., Aumüller et al., 2020) for the top-performing ANN methods. The inclusion criteria are: (i) the method is published as a library with Python bindings; (ii) the library is open source; (iii) the method has top performance in its class or is widely used in applications. See Table 3 for the included methods.

Hyperparameter settings. Most implementations contain several hyperparameters that have to be tuned to a specific workload, and VIBE allows for easy configuration of a grid search for these hyperparameters. Initial choices for applicable methods were obtained from the hyperparameter grids of the ANN-Benchmarks project (Aumüller et al., 2020), and for new methods from the library documentation and the respective research papers. We then expanded these grids as necessary based on intermediate results. To attest to the robustness of an implementation, the hyperparameters are shared among all datasets. The hyperparameter choices are available in our code repository. In the figures of this section, we report the Pareto frontiers over all hyperparameter choices.

5.2 In-distribution setting

We compare the ANN algorithms on 12 in-distribution datasets. An overview of the results (on 6 representative datasets) can be found in Figure 1. For selected algorithms, the recall-QPS curves are reported in Figure 2 for two representative text embedding datasets and in Figure 3 for two representative image embedding datasets. The results for all datasets can be found in Appendix F.1. For results on all algorithms, see our companion website or Appendix F.2 for results on two datasets.

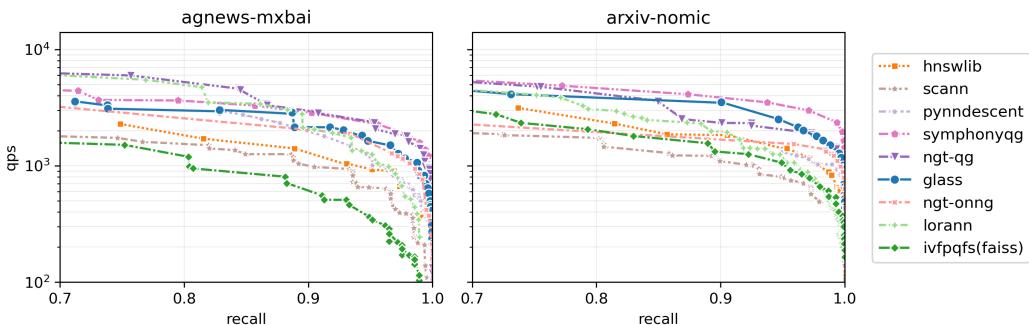


Figure 2: Recall/throughput tradeoff on two **text embedding** datasets (in-distribution queries). Graph-based SymphonyQG, Glass, and NGT-QG have the highest throughput at average recall above 90%.

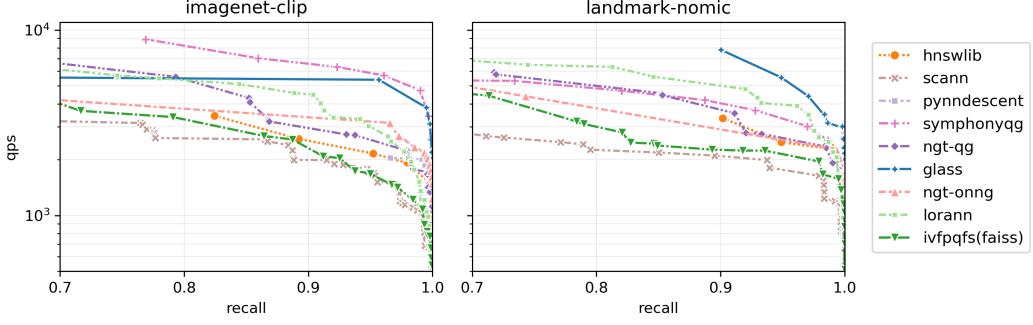


Figure 3: Recall/throughput tradeoff on two **image embedding** datasets (in-distribution queries). Graph-based SymphonyQG and Glass, and clustering-based LoRANN have the highest throughput.

In line with the earlier performance evaluations (e.g., Li et al., 2019; Aumüller et al., 2020), the general trend is that the graph and clustering methods significantly outperform the hashing and tree methods: the top-performing methods contain only graph methods and clustering methods⁵. In particular, the graph methods SymphonyQG and Glass have the top performance: at 95% recall, SymphonyQG has the highest throughput on 5 and Glass on 4 out of 12 datasets (graph-based NGT-QG is the fastest on 2 datasets, and clustering-based LoRANN on 1 dataset). However, as shown in Appendix E, graph-based approaches have the disadvantage of longer index build times.

Robustness of algorithms. The performance difference between the 100 hardest and the 100 easiest queries, as measured by relative contrast (RC), is shown in Figure 4 for the four top-performing algorithms at 90% recall. Glass and LoRANN return almost correct results for the easiest queries, but have just over 70% recall for the hardest queries. In contrast, SymphonyQG and NGT-QG have more robust performance: the RC of queries does not seem to affect the recall of SymphonyQG⁶.

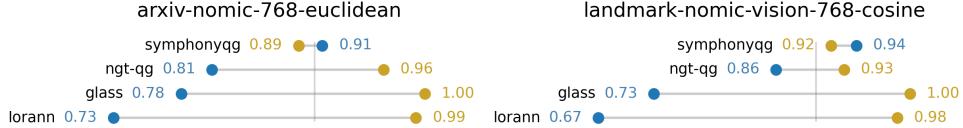


Figure 4: Recall of the 100 hardest and easiest queries (by RC score) on two datasets, considering the fastest configuration achieving an average recall of at least 90% (marked by the vertical line). SymphonyQG is robust w.r.t. the query difficulty, while Glass and LoRANN show larger variability.

Higher throughput with binary datasets and GPU deployment. VIBE also includes support for binary datasets (with Hamming distance) and GPU algorithms (see Section 4.2) that can be used to achieve a higher throughput: for example, on the float-precision agnews-mxbai dataset at 90% recall, the best algorithm achieves a throughput of 2×10^3 QPS (Figure 2), whereas on the binarized version of the same dataset, the best algorithm achieves a throughput of 8×10^3 QPS (Figure 5, left panel). Using an NVIDIA V100 GPU, the best GPU algorithm achieves a throughput of 10^5 QPS when all 1000 test queries are given as a batch (Figure 5, right panel). More results are shown in Appendix F.3: on binary datasets, the graph method NGT-ONNG (Iwasaki and Miyazaki, 2018) is the fastest algorithm, while the graph-based CAGRA (Ootomo et al., 2024) is the fastest GPU method.

5.3 Out-of-distribution setting

We compare the performance of ANN algorithms on 6 out-of-distribution (OOD) datasets. An overview of the results at 95% recall can be found in Figure 1. The results for two text-to-image datasets are shown in Figure 7 and for the approximate attention computation datasets in Figure 8.

Text-to-image. On the text-to-image OOD datasets, the most efficient regular ANN methods outperform the ANN methods specifically tailored for OOD data (see Figure 7). However, as can be

⁵For clarity, in the figures of this section, we report the results only for selected 9 methods. The complete results can be found on the website and (for a couple of representative datasets) in Appendix F.2.

⁶SymphonyQG's better performance on the harder queries is likely due to random variation.

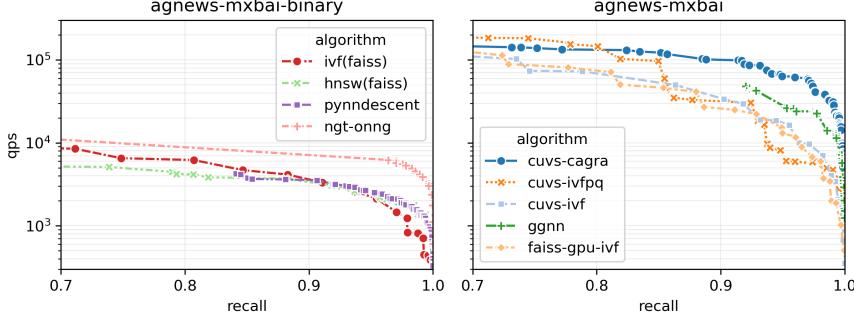


Figure 5: Left: recall/throughput tradeoff on **binary data**. Right: recall/throughput tradeoff of **GPU** algorithms. The graph-based NGT-ONNG has the highest throughput on the binary data, and the graph-based CAGRA is the fastest method in the GPU setting.

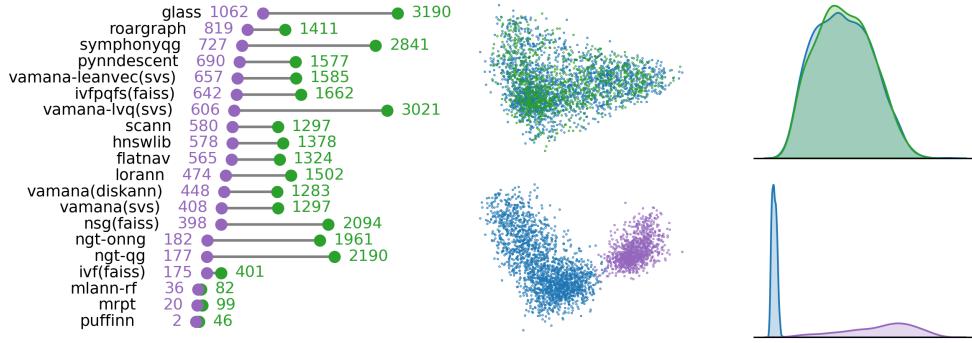


Figure 6: Difference in throughput (QPS) between out-of-distribution and in-distribution queries on the laion-clip dataset at 95% recall. The middle column plots show the corpus and the queries (both in-distribution and out-of-distribution) projected into two dimensions using PCA, while the rightmost plots report the distribution of Mahalanobis distances of between the corpus points (blue) and the queries and the corpus points (green and purple). The performance of the top regular ANN methods, such as Glass, SymphonyQG and NGT-QG degrades significantly on the OOD queries.

seen from Figure 6, the performance of the regular ANN algorithms degrades on out-of-distribution queries compared to in-distribution queries (for the results on additional datasets, see Appendix G). This performance gap suggests that there is still a large room for improvement in the development of out-of-distribution ANN algorithms.

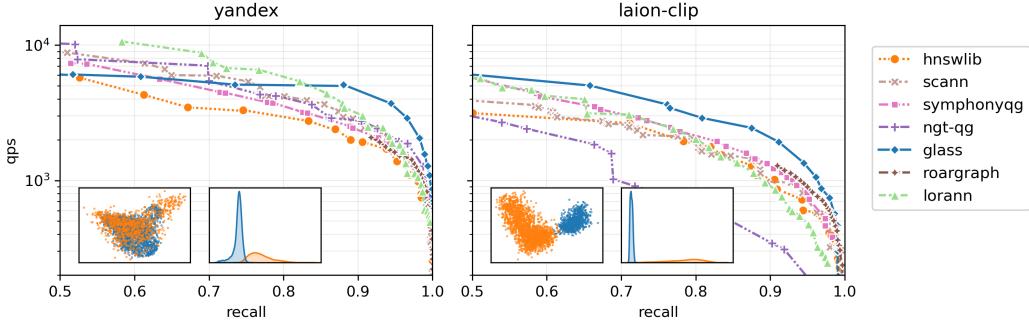


Figure 7: Recall/throughput tradeoff on two **text-to-image** datasets, with **out-of-distribution** queries. The leftmost insets show PCA projections of the corpus points and the queries, and the rightmost insets show the distributions of the Mahalanobis distances between the corpus points (blue), and the queries and the corpus points (orange). At high recall levels, Glass has the highest throughput on both datasets, outperforming the specialized OOD methods, such as RoarGraph.

Approximate attention computation. Regular ANN methods struggle on the approximate attention computation datasets (see Figure 8): for example, graph methods Glass and hnswlib do not even reach 50% recall with reasonable hyperparameter choices. Clustering-based methods perform relatively better, but ScaNN and IVF-PQ perform worse than regular IVF. In contrast, the methods that are specifically designed for OOD data, i.e., RoarGraph, MLANN-RF, and LoRANN, reach the highest recall levels and have the top performance. The performance difference between the OOD methods and the regular ANN methods is more pronounced on the llama data set where the difference between the corpus distribution and the query distribution is larger (Figure 8, right panel), compared to the yi data set where the difference is smaller (Figure 8, left panel).

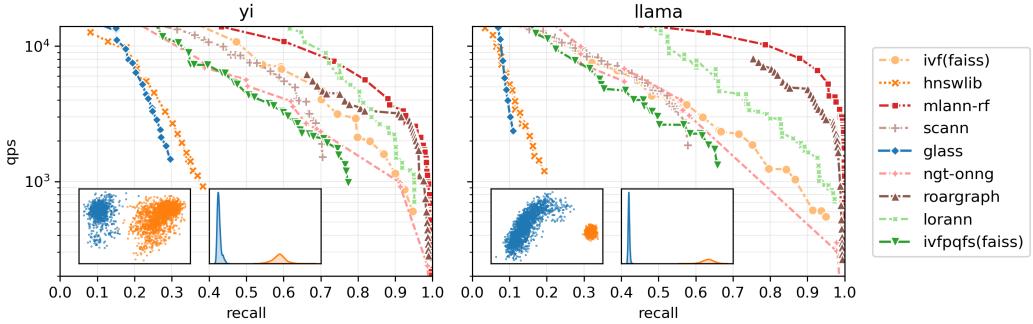


Figure 8: Recall/throughput tradeoff on the **approximate attention computation** datasets, with **out-of-distribution** queries. The specialized OOD methods MLANN-RF, RoarGraph, and LoRANN outperform the regular ANN methods, such as Glass, that struggle to reach the highest recall levels.

6 Discussion

6.1 Key findings

Comparison of different classes of methods. The graph-based SymphonyQG (Gou et al., 2025) and Glass (Wang, 2025) are the top-performing implementations. SymphonyQG is also the most robust method w.r.t. the query difficulty (see Figure 4). Graph and clustering indexes outperform both tree and hashing indexes; this is in line with earlier research (Aumüller et al., 2020; Li et al., 2019).

Reduced precision. All top graph (SymphonyQG, Glass, NGT-QG) and clustering (LoRANN, ScaNN, Faiss-IVF-PQ) methods use different types of quantization in at least some of their computations. SymphonyQG uses RaBitQ (Gao and Long, 2024), while NGT-QG, ScaNN, and Faiss-IVF-PQ employ efficient 4-bit PQ implementations (André et al., 2017). LoRANN uses 8-bit scalar quantization to speed up vector-matrix multiplications, and Glass, an efficient HNSW implementation using scalar quantization, significantly outperforms the full-precision HNSW implementations.

OOD setting. In the out-of-distribution (OOD) setting, two trends emerge: (1) On the approximate attention computation use case, where the difference between the corpus and query distributions is the largest, regular ANN algorithms struggle to reach the highest recall levels, and the algorithms specifically tailored for OOD setting outperform them. (2) However, on the other OOD datasets, our results suggest that while the performance of the regular ANN algorithms degrades in the OOD setting, the current OOD algorithms do not yield significant performance improvements over them.

6.2 Limitations

Specialized settings. We limit the scope of the benchmark to standard in-distribution and out-of-distribution ANN search and do not consider specialized settings, such as filtered, streaming, and sparse search (see, e.g., Simhadri et al., 2024).

Large datasets. Due to the high dimensionality of embeddings, we focus on million-scale datasets to keep the computational requirements of the benchmark reasonable and to allow us to study the behavior of a large set of algorithms on a wide array of embeddings. However, evaluating the performance of ANN algorithms on billion-scale datasets to see whether the key findings of this article also hold in that regime is an interesting direction for future research.

Acknowledgments

This work has been supported by the Research Council of Finland (grant #361902 and the Flagship programme: Finnish Center for Artificial Intelligence FCAI) and the Jane and Aatos Erkko Foundation (BioDesign project, grant #7001702). Martin Aumüller received funding from the Innovation Fund Denmark for the project DIREC (9142-00001B). The authors wish to acknowledge CSC – IT Center for Science, Finland, for computational resources.

References

- Cecilia Aguerrebere, Mark Hildebrand, Ishwar Singh Bhati, Theodore Willke, and Mariano Tepper. Locally-adaptive quantization for streaming vector search. *arXiv preprint arXiv:2402.02044*, 2024.
- Fabien André, Anne-Marie Kermarrec, and Nicolas Le Scouarnec. Accelerated nearest neighbor search with quick ADC. In *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval*, pages 159–166, 2017.
- Martin Aumüller and Matteo Ceccarello. The role of local dimensionality measures in benchmarking nearest neighbor search. *Information Systems*, 101:101807, 2021.
- Martin Aumüller and Matteo Ceccarello. Recent approaches and trends in approximate nearest neighbor search, with remarks on benchmarking. *IEEE Data Engineering Bulletin*, 47(3):89–105, 2023.
- Martin Aumüller, Tobias Christiani, Rasmus Pagh, and Michael Vesterli. PUFFINN: Parameterless and universally fast finding of nearest neighbors. In *27th Annual European Symposium on Algorithms (ESA 2019)*, pages 10:1–10:16. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2019.
- Martin Aumüller, Erik Bernhardsson, and Alexander Faithfull. ANN-Benchmarks: A benchmarking tool for approximate nearest neighbor algorithms. *Information Systems*, 87:101374, 2020.
- Ilias Azizi, Karima Echihabi, and Themis Palpanas. Graph-based vector search: An experimental evaluation of the state-of-the-art. *Proceedings of the ACM on Management of Data*, 3(1):1–31, 2025.
- Amanda Bertsch, Uri Alon, Graham Neubig, and Matthew Gormley. Unlimiformer: Long-range transformers with unlimited length input. *Advances in Neural Information Processing Systems*, 36: 35522–35543, 2023.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. In *International Conference on Machine Learning*, pages 2206–2240. PMLR, 2022.
- Lawrence Cayton and Sanjoy Dasgupta. A learning framework for nearest neighbor search. *Advances in Neural Information Processing Systems*, 20:233–240, 2007.
- Meng Chen, Kai Zhang, Zhenying He, Yinan Jing, and X. Sean Wang. RoarGraph: A projected bipartite graph for efficient cross-modal approximate nearest neighbor search. *Proceedings of the VLDB Endowment*, 17(11):2735–2749, 2024.
- Paul Covington, Jay Adams, and Emre Sargin. Deep neural networks for YouTube recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*, pages 191–198, 2016.
- Sanjoy Dasgupta and Yoav Freund. Random projection trees and low dimensional manifolds. In *Proceedings of the Fortieth Annual ACM Symposium on Theory of Computing*, pages 537–546, 2008.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. IEEE, 2009.

Wei Dong, Charikar Moses, and Kai Li. Efficient k-nearest neighbor graph construction for generic similarity measures. In *Proceedings of the 20th International Conference on World Wide Web*, page 577–586, 2011.

Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvassy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. *arXiv preprint arXiv:2401.08281*, 2024.

Alex Fang, Gabriel Ilharco, Mitchell Wortsman, Yuhao Wan, Vaishaal Shankar, Achal Dave, and Ludwig Schmidt. Data determines distributional robustness in contrastive language image pre-training (CLIP). In *International Conference on Machine Learning*, pages 6216–6234. PMLR, 2022.

Jerome H. Friedman, Jon Louis Bentley, and Raphael Ari Finkel. An algorithm for finding best matches in logarithmic expected time. *ACM Transactions on Mathematical Software*, 3(3):209–226, 1977.

Cong Fu, Chao Xiang, Changxu Wang, and Deng Cai. Fast approximate nearest neighbor search with the navigating spreading-out graphs. *PVLDB*, 12(5):461 – 474, 2019.

Jianyang Gao and Cheng Long. RaBitQ: Quantizing high-dimensional vectors with a theoretical error bound for approximate nearest neighbor search. *Proceedings of the ACM on Management of Data*, 2(3):1–27, 2024.

Yutong Gou, Jianyang Gao, Yuexuan Xu, and Cheng Long. SymphonyQG: Towards symphonious integration of quantization and graph for approximate nearest neighbor search. *Proceedings of the ACM on Management of Data*, 3(1):1–26, 2025.

Fabian Groh, Lukas Ruppert, Patrick Wieschollek, and Hendrik P.A. Lensch. GGNN: Graph-based GPU nearest neighbor search. *IEEE Transactions on Big Data*, 9(1):267–279, 2022.

Michael Günther, Jackmin Ong, Isabelle Mohr, Alaaeddine Abdessalem, Tanguy Abel, Mohammad Kalim Akram, Susana Guzman, Georgios Mastrapas, Saba Sturua, Bo Wang, et al. Jina embeddings 2: 8192-token general-purpose text embeddings for long documents. *arXiv preprint arXiv:2310.19923*, 2023.

Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. Accelerating large-scale inference with anisotropic vector quantization. In *International Conference on Machine Learning*, pages 3887–3896. PMLR, 2020.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented language model pre-training. In *International Conference on Machine Learning*, pages 3929–3938. PMLR, 2020.

Junfeng He, Sanjiv Kumar, and Shih-Fu Chang. On the difficulty of nearest neighbor search. In *International Conference on Machine Learning*. PMLR, 2012.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.

Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. CodeSearchNet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436*, 2019.

Ville Hyvönen, Teemu Pitkänen, Sotiris Tasoulis, Elias Jääsaari, Risto Tuomainen, Liang Wang, Jukka Corander, and Teemu Roos. Fast nearest neighbor search through sparse random projections and voting. In *Proceedings of the 2016 IEEE International Conference on Big Data*, pages 881–888. IEEE, 2016.

Ville Hyvönen, Elias Jääsaari, and Teemu Roos. A multilabel classification framework for approximate nearest neighbor search. *Advances in Neural Information Processing Systems*, 35: 35741–35754, 2022.

- Masajiro Iwasaki and Daisuke Miyazaki. Optimization of indexing based on k-nearest neighbor graph for proximity search in high-dimensional data. *arXiv preprint arXiv:1810.07355*, 2018.
- Elias Jääsaari, Ville Hyvönen, and Teemu Roos. LoRANN: Low-rank matrix factorization for approximate nearest neighbor search. *Advances in Neural Information Processing Systems*, 37: 102121–102153, 2024.
- Shikhar Jaiswal, Ravishankar Krishnaswamy, Ankit Garg, Harsha Vardhan Simhadri, and Sheshansh Agrawal. OOD-DiskANN: Efficient and scalable graph ANNS for out-of-distribution queries. *arXiv preprint arXiv:2211.12850*, 2022.
- Herve Jegou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(1):117–128, 2011.
- Chao Jia, Yinfai Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.
- Daniel Khashabi, Amos Ng, Tushar Khot, Ashish Sabharwal, Hannaneh Hajishirzi, and Chris Callison-Burch. GooAQ: Open question answering with diverse answer types. *arXiv preprint arXiv:2104.08727*, 2021.
- Sean Lee, Aamir Shakir, Darius Koenig, and Julius Lipp. Open source strikes bread - new fluffy embeddings model, 2024. URL <https://www.mixedbread.ai/blog/mxbai-embed-large-v1>.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474, 2020.
- Patrick Lewis, Yuxiang Wu, Linqing Liu, Pasquale Minervini, Heinrich Küttler, Aleksandra Piktus, Pontus Stenetorp, and Sebastian Riedel. PAQ: 65 million probably-asked questions and what you can do with them. *Transactions of the Association for Computational Linguistics*, 9:1098–1115, 2021.
- Minghan Li, Xilun Chen, Ari Holtzman, Beidi Chen, Jimmy Lin, Scott Yih, and Victoria Lin. Nearest neighbor speculative decoding for LLM generation and attribution. *Advances in Neural Information Processing Systems*, 38:80987–81015, 2024.
- Wen Li, Ying Zhang, Yifang Sun, Wei Wang, Mingjie Li, Wenjie Zhang, and Xuemin Lin. Approximate nearest neighbor search on high dimensional data—experiments, analyses, and improvement. *IEEE Transactions on Knowledge and Data Engineering*, 32(8):1475–1488, 2019.
- Xianming Li and Jing Li. Angle-optimized text embeddings. *arXiv preprint arXiv:2309.12871*, 2023.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference*, pages 740–755. Springer, 2014.
- Di Liu, Meng Chen, Baotong Lu, Huiqiang Jiang, Zhenhua Han, Qianxi Zhang, Qi Chen, Chenguodong Zhang, Bailu Ding, Kai Zhang, et al. RetrievalAttention: Accelerating long-context LLM inference via vector retrieval. *arXiv preprint arXiv:2409.10516*, 2024.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- Alexander Long, Wei Yin, Thalaiyasingam Ajanthan, Vu Nguyen, Pulak Purkait, Ravi Garg, Alan Blair, Chunhua Shen, and Anton van den Hengel. Retrieval augmented classification for long-tail visual recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6959–6969, 2022.

- Yu A. Malkov and Dmitry A. Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4):824–836, 2018.
- Pierre-Emmanuel Mazaré, Gergely Szilvasy, Maria Lomeli, Francisco Massa, Naila Murray, Hervé Jégou, and Matthijs Douze. Inference-time sparse attention with asymmetric indexing. *arXiv preprint arXiv:2502.08246*, 2025.
- Leland McInnes. PyNNDescent: A python library for approximate nearest neighbors, 2018. URL <https://github.com/lmcinnes/pynndescent/>.
- Blaise Munyampirwa, Vihan Lakshman, and Benjamin Coleman. Down with the hierarchy: The 'H' in HNSW stands for "Hubs". *arXiv preprint arXiv:2412.01940*, 2024.
- Zach Nussbaum, Brandon Duderstadt, and Andriy Mulyar. Nomic embed vision: Expanding the latent space. *arXiv preprint arXiv:2406.18587*, 2024.
- Zach Nussbaum, John X Morris, Brandon Duderstadt, and Andriy Mulyar. Nomic embed: Training a reproducible long context text embedder. *Transactions of Machine Learning Research*, 2025.
- NVIDIA. cuVS: Vector search and clustering on the GPU, 2024. URL <https://github.com/rapidsai/cuvs>.
- Hiroyuki Ootomo, Akira Naruse, Corey Nolet, Ray Wang, Tamas Feher, and Yong Wang. CAGRA: Highly parallel graph construction and approximate nearest neighbor search for GPUs. In *2024 IEEE 40th International Conference on Data Engineering (ICDE)*, pages 4236–4247. IEEE, 2024.
- OpenAI. text-embedding-3-large, 2024. URL <https://platform.openai.com/docs/guides/embeddings>. OpenAI Embeddings Model.
- Maxime Oquab, Timothée Darct, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Noubi, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOV2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
- Ninh Pham and Tao Liu. Falconn++: A locality-sensitive filtering approach for approximate nearest neighbor search. *Advances in Neural Information Processing Systems*, 35:31186–31198, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- Minjoon Seo, Jinhyuk Lee, Tom Kwiatkowski, Ankur Parikh, Ali Farhadi, and Hannaneh Hajishirzi. Real-time open-domain question answering with dense-sparse phrase index. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2019.
- Harsha Vardhan Simhadri, George Williams, Martin Aumüller, Matthijs Douze, Artem Babenko, Dmitry Baranchuk, Qi Chen, Lucas Hosseini, Ravishankar Krishnaswamy, Gopal Srinivasa, et al. Results of the NeurIPS’21 challenge on billion-scale approximate nearest neighbor search. In *NeurIPS 2021 Competitions and Demonstrations Track*, pages 177–189. PMLR, 2022.
- Harsha Vardhan Simhadri, Martin Aumüller, Amir Ingber, Matthijs Douze, George Williams, Magdalene Dobson Manohar, Dmitry Baranchuk, Edo Liberty, Frank Liu, Ben Landrum, et al. Results of the Big ANN: NeurIPS’23 competition. *arXiv preprint arXiv:2409.17424*, 2024.

- Spotify. Annoy: Approximate nearest neighbors oh yeah, 2013. URL <https://github.com/spotify/annoy>.
- Robert F. Sproull. Refinements to nearest-neighbor searching in k-dimensional trees. *Algorithmica*, 6:579–589, 1991.
- Suhas Jayaram Subramanya, Fnu Devvrit, Harsha Vardhan Simhadri, Ravishankar Krishnawamy, and Rohan Kadekodi. DiskANN: Fast accurate billion-point nearest neighbor search on a single node. *Advances in Neural Information Processing Systems*, 32, 2019.
- Mariano Tepper, Ishwar Singh Bhati, Cecilia Aguerrebere, Mark Hildebrand, and Theodore L. Willke. LeanVec: Searching vectors faster by making them fit. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- Boxin Wang, Wei Ping, Lawrence McAfee, Peng Xu, Bo Li, Mohammad Shoeybi, and Bryan Catanzaro. InstructRetro: Instruction tuning post retrieval-augmented pretraining. In *International Conference on Machine Learning*, pages 51255–51272. PMLR, 2024.
- Mengzhao Wang, Xiaoliang Xu, Qiang Yue, and Yuxiang Wang. A comprehensive survey and experimental comparison of graph-based approximate nearest neighbor search. *Proceedings of the VLDB Endowment*, 14(11):1964–1978, 2021.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. MiniLM: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in Neural Information Processing Systems*, 33:5776–5788, 2020.
- Zihao Wang. Graph library for approximate similarity search, 2025. URL <https://github.com/zilliztech/pyglass>.
- Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2 - a large-scale benchmark for instance-level recognition and retrieval. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2575–2584, 2020.
- Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N Bennett, Junaid Ahmed, and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning for dense text retrieval. In *International Conference on Learning Representations*, 2021.
- Yahoo. NGT: Neighborhood graph and tree for indexing high-dimensional data, 2015. URL <https://github.com/yahoojapan/NGT/>.
- Ji Yang, Xinyang Yi, Derek Zhiyuan Cheng, Lichan Hong, Yang Li, Simon Xiaoming Wang, Taibai Xu, and Ed H Chi. Mixed negative sampling for learning two-tower neural networks in recommendations. In *Companion Proceedings of the Web Conference 2020*, pages 441–447, 2020.
- Amir Zandieh, Insu Han, Majid Daliri, and Amin Karbasi. KDEformer: Accelerating transformers via kernel density estimation. In *International Conference on Machine Learning*, pages 40605–40623. PMLR, 2023.
- Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. A robustly optimized BERT pre-training approach with post-training. In *Proceedings of the 20th Chinese National Conference on Computational Linguistics*, pages 1218–1227, 2021.

A Interactive website

Our interactive website is available at <https://vector-index-bench.github.io/>. It sports several pages that allow to:

- compare all the algorithms at a glance at a given configurable recall level;
- explore the distribution of queries and data of different datasets;
- investigate the tradeoffs of different implementations (Figure 9 shows a screenshot);
- assess the performance of all the configurations we tested for each algorithm (see Figure 10);
- evaluate the robustness of all the algorithms on easy and difficult queries.

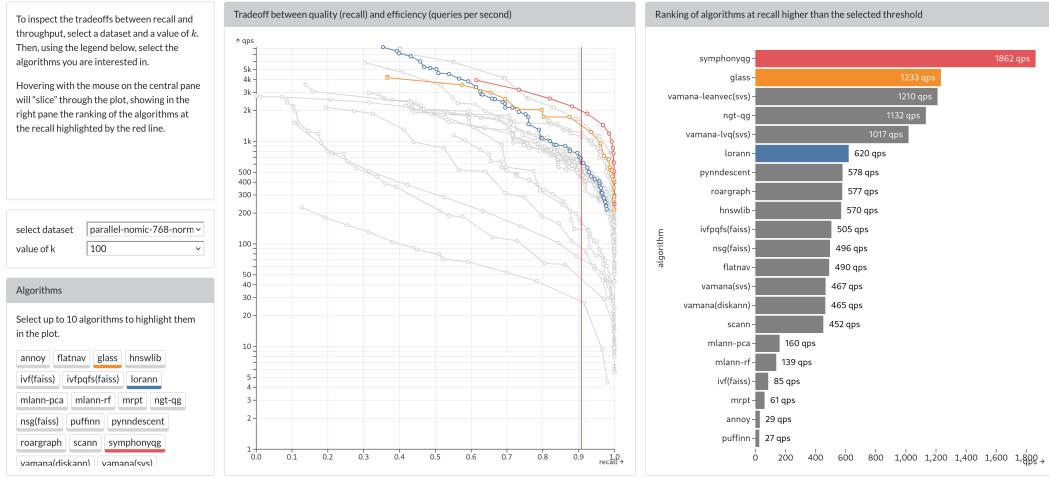


Figure 9: Screenshot of the webpage that allows to interactively explore the recall/throughput tradeoff. The controls on the left allow selecting the dataset, the value of k , and the algorithms to highlight. The central plot reports the tradeoff curves. Hovering with the mouse on the central plot moves the vertical red line which marks a recall threshold: the bar chart on the right is then dynamically updated to report the fastest configuration of each algorithm attaining at least the selected recall.

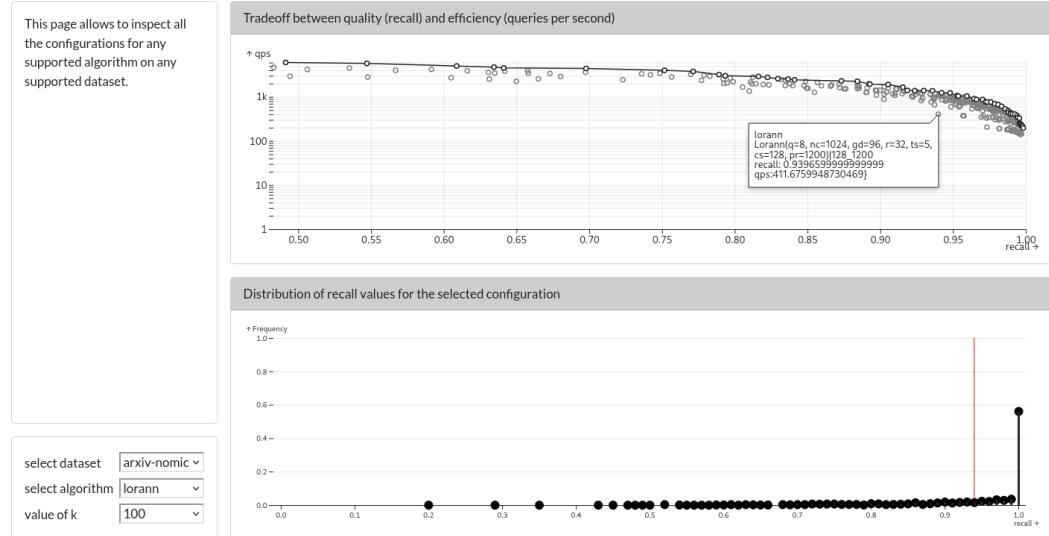


Figure 10: Screenshot of the webpage that allows to interactively explore the behavior of all the tested configurations of a given algorithm. Each dot represents a different configuration, which is revealed when hovering with the mouse over it. At the same time, the bottom plot reports the distribution of the recall values attained by each individual query with the selected configuration.

B Dataset information

In this section, we detail the embedding models and data sources used to produce the VIBE datasets. Most of the datasets in VIBE are created by us using the embeddings and data sources detailed in this section, and our embedding outputs are distributed under CC BY 4.0. However, as part of VIBE we also redistribute the following embedding datasets:

- GloVe (Pennington et al., 2014) (under PDDL 1.0)
- Yandex Text-to-Image (Simhadri et al., 2024) (a subset of 1 million vectors, under CC BY 4.0)
- LAION (a subset of 1 million vectors, under CC BY 4.0)

The embedding models used to produce the new datasets (see Tables 1 and 2) in VIBE are shown in Table 4. For a brief discussion, we refer to Section 4.1.

Table 4: Embedding models used to produce the VIBE embedding datasets. The license for each embedding applies only to the model itself; none of the models claim rights over their outputs.

Name	Model	Citation	License
MXBAI	mxbai-embed-large-v1	Li and Li (2023)	Apache 2.0
Nomic Text	nomic-embed-text-v1.5	Nussbaum et al. (2025)	Apache 2.0
DistilRoBERTa	all-distilroberta-v1	Sanh et al. (2019)	Apache 2.0
CLIP	clip-vit-base-patch32	Radford et al. (2021)	MIT
Nomic Vision	nomic-embed-vision-v1.5	Nussbaum et al. (2024)	Apache 2.0
MiniLM	all-MiniLM-L6-v2	Wang et al. (2020)	Apache 2.0
ResNet	resnet50	He et al. (2016)	BSD 3
Jina	jina-embeddings-v2-base-code	Günther et al. (2023)	Apache 2.0
DINOv2	vit_base_patch16_224.dino	Oquab et al. (2024)	Apache 2.0
OpenAI	text-embedding-3-large	OpenAI (2024)	Custom
ALIGN	align-base	Jia et al. (2021)	Custom

Table 5: Data sources used to produce the VIBE embedding datasets.

Name	Type	<i>n</i>	Reference	Citation
AGNews	Text	769 382	¹	N/A
arXiv abstracts	Text	1 344 643	²	N/A
CCNews	Text	495 328	³	N/A
CelebA	Image	201 599	⁴	Liu et al. (2015)
CodeSearchNet	Code	1 374 067	⁵	Husain et al. (2019)
Google Q&A	Text	1 475 024	⁶	Khashabi et al. (2021)
ImageNet	Image	1 281 167	⁷	Deng et al. (2009)
ImageNet-Captions	Text	316 648	⁸	Fang et al. (2022)
Landmark	Image	760 757	⁹	Weyand et al. (2020)
MS COCO	Text/Image	282 360	¹⁰	Lin et al. (2014)
Simple Wiki	Text	260 372	¹¹	N/A
Yahoo Q&A	Text	677 305	¹²	N/A

¹ http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html

² <https://www.kaggle.com/datasets/Cornell-University/arxiv>

³ <https://commoncrawl.org/blog/news-dataset-available>

⁴ <https://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

⁵ <https://github.com/github/CodeSearchNet>

⁶ <https://github.com/allenai/gooaq>

⁷ <https://www.image-net.org/>

⁸ <https://github.com/mlfoundations/imagenet-captions>

⁹ <https://github.com/cvdfoundation/google-landmark>

¹⁰ <https://cocodataset.org/>

¹¹ <https://dumps.wikimedia.org/>

¹² <https://www.kaggle.com/datasets/soumikrakshit/yahoo-answers-dataset>

The used data sources are shown in Table 5. Most of the data sources are under custom non-commercial licenses; for details, we refer to the provided references for each data source.

For our approximate attention computation datasets, we use two long-context LLMs:

- **yi-128-ip**: Extracted from Layer 32, Head 14 of the 01-ai/Yi-6B-200K model using a copy of the book *Pride and Prejudice*⁷ as the prompt.
- **llama-128-ip**: Extracted from Layer 13, Head 16 of the gradientai/Llama-3-8B-Instruct-262k model using a copy of the book *The Picture of Dorian Gray*⁸ as the prompt.

Following Mazaré et al. (2025), we discard the keys and queries corresponding to the first token and the last 2047 tokens, as typically they have higher inner products with the queries and should be computed exactly.

⁷<https://www.gutenberg.org/ebooks/1342>

⁸<https://www.gutenberg.org/ebooks/174>

C Dataset difficulty

Figure 11 summarizes the distribution of relative contrast values (cf. Section 2.1) over the different datasets. Visually, datasets stemming from image embeddings have a wider range of relative contrast values than text embeddings. As demonstrated in Section 5.2, for a single dataset the relative contrast values of query points are a good indicator of the difficulty of answering the query accurately.

To provide further context on the comparability between datasets we focus on four datasets `arxiv-nomic`, `gooaq-distilroberta`, `codesearchnet-jina`, and `glove` and three implementations SymphonyQG, LoRANN, and ScaNN. The first three datasets are similar in the number of vectors and their dimensionality, while `glove` has both lower dimensionality and contains fewer vectors. However, from the distribution of RC values, queries of `glove` have the lowest contrast. The experiments (see the full set of figures in Appendix F) confirm the difficulty of the dataset. At average recall of at least 90%, the highest throughput of ScaNN is 1049 QPS, SymphonyQG achieves 560 QPS, and no hyperparameter setting in LoRANN achieves the necessary quality. In contrast, SymphonyQG achieves at least 3200 QPS on the other three datasets, ScaNN achieves roughly the same throughput, while LoRANN achieves more than 1600 QPS on `arxiv` and `gooaq-distilroberta`, and 885 QPS on `codesearchnet-jina`.

As noticed by Aumüller and Ceccarello (2021), comparing relative contrast values between different datasets can lead to misleading predictions: While `imagenet-clip` is indeed easier than its out-of-distribution companion `imagenet-align` (SymphonyQG achieves a maximum throughput of 6320 QPS on the former, and 1805 QPS on the latter), the throughput achieved is higher than the measurements observed for `glove`, despite the relative contrast predicting less contrast for `imagenet-align`.

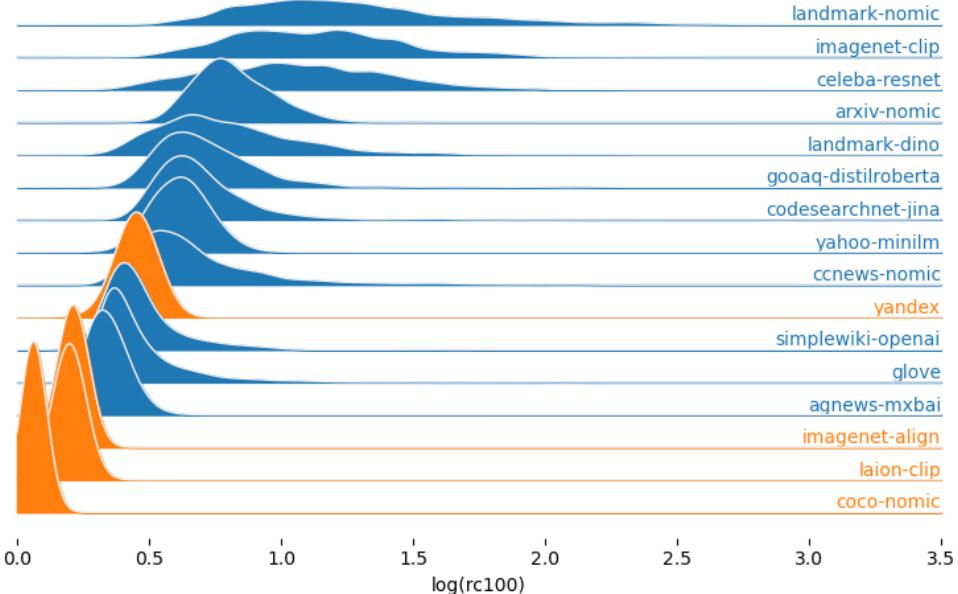


Figure 11: Distribution of relative contrast @ 100 for in-distribution and out-of-distribution datasets. Density curves are arranged by the median relative contrast.

D Algorithms

In Table 6, we give a detailed list of versions of each algorithm used in the experiments. Clustering-based methods are highlighted in green, tree-based in purple, hashing-based in orange, and graph-based in blue.

Table 6: Details of the algorithm implementations included in the benchmark.

Method	Reference	Version
ANNOY	Spotify (2013)	1.17.3
FALCONN++	Pham and Liu (2022)	git+5fd3f17
FlatNav	Munyampirwa et al. (2024)	0.1.2
Glass	Wang (2025)	1.0.5
HNSW	Malkov and Yashunin (2018)	0.8.0
IVF	Jegou et al. (2011); Douze et al. (2024)	1.11.0
IVF-PQ	Jegou et al. (2011); Douze et al. (2024)	1.11.0
LVQ	Aguerrebere et al. (2024)	0.0.7
LeanVec	Tepper et al. (2024)	0.0.7
LoRANN	Jääsaari et al. (2024)	0.2
MLANN	Hyvönen et al. (2022)	git+40848e7
MRPT	Hyvönen et al. (2016)	2.0.1
NGT-ONNG	Iwasaki and Miyazaki (2018); Yahoo (2015)	git+83d5896
NGT-QG	Yahoo (2015)	git+83d5896
NSG	Fu et al. (2019); Douze et al. (2024)	1.11.0
PUFFINN	Aumüller et al. (2019)	git+fd86b0d
PyNNDescent	Dong et al. (2011); McInnes (2018)	0.5.13
RoarGraph	Chen et al. (2024)	git+f2b49b6
ScaNN	Guo et al. (2020)	1.4.0
SymphonyQG	Gou et al. (2025)	git+32a0019
Vamana	Subramanya et al. (2019)	0.7.0

The list of algorithms used in the GPU experiments is given in Table 7.

Table 7: Details of the algorithm implementations included in the benchmark.

Method	Reference	Version
CAGRA	Ootomo et al. (2024); NVIDIA (2024)	25.4.0
GGNN	Groh et al. (2022)	0.9
IVF-PQ (cuVS)	NVIDIA (2024)	25.4.0
IVF (cuVS)	NVIDIA (2024)	25.4.0
IVF (Faiss)	Johnson et al. (2019); Douze et al. (2024)	1.11.0

E Index construction time

We measure the index construction time on a single CPU core for all methods. Table 8 shows index construction times for all algorithms on six datasets at 90% recall for $k = 100$ at the optimal hyperparameters (a dash means that the algorithm did not reach the necessary recall).

Clustering-based methods are highlighted in green, tree-based in purple, hashing-based in orange, and graph-based in blue.

Table 8: Index construction times (seconds) with throughput-optimized hyperparameters at 90% recall. Algorithms are sorted based on their average index construction times.

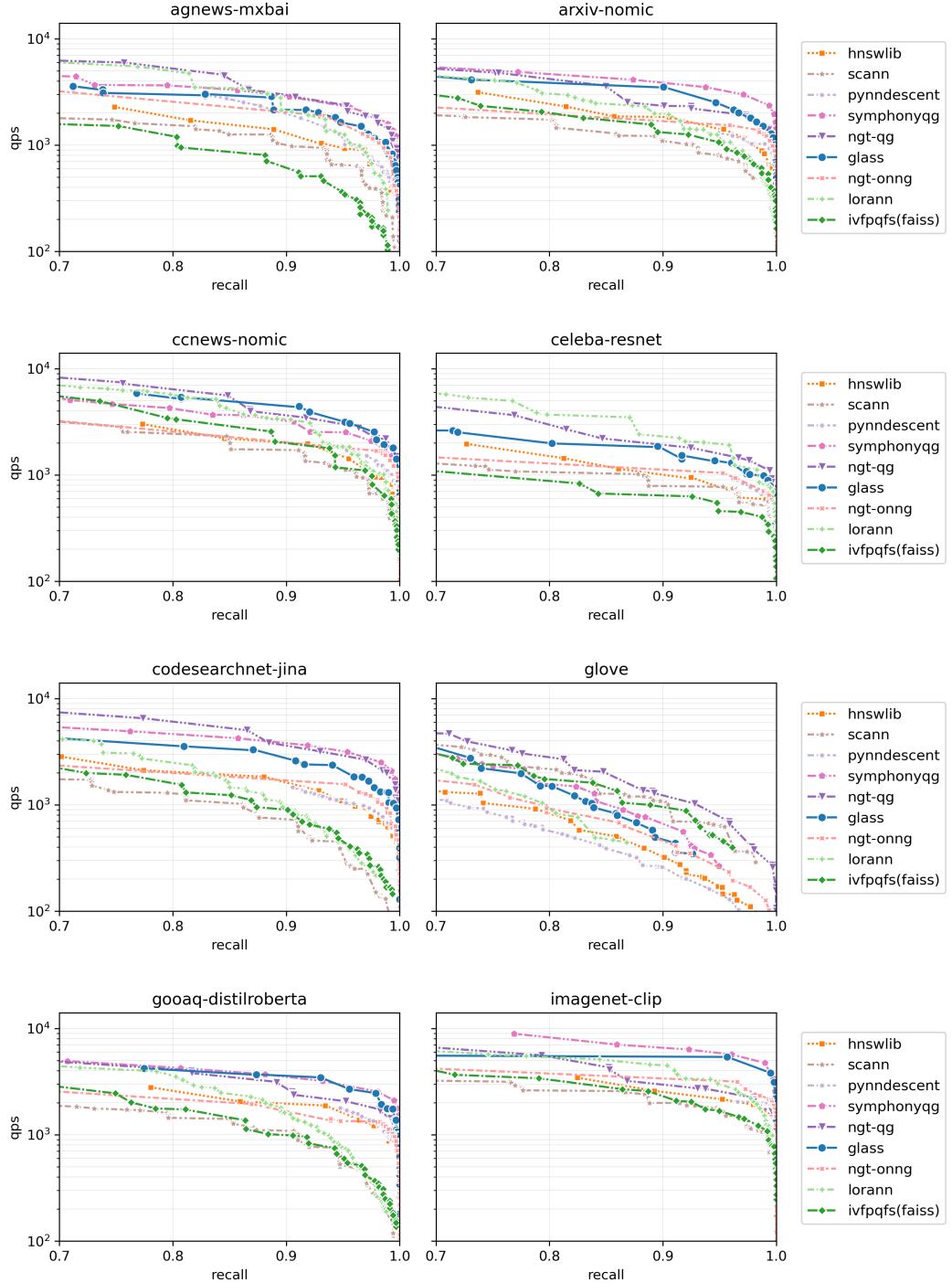
algorithm	agnews	arxiv	gooaq	imagenet	landmark	yahoo
ivfpqfs(faiss)	64	68	121	45	126	19
ivf(faiss)	84	214	229	156	203	89
annoy	314	1007	1136	292	1253	280
scann	838	984	1125	626	556	256
vamana-leanvec(svs)	456	1385	1190	1102	440	479
lorann	529	1057	1262	1096	469	678
glass	1001	1377	1532	751	334	813
pynndescent	361	1506	1753	1185	818	516
vamana-lvq(svs)	1570	2462	2248	703	519	705
vamana(diskann)	1570	2769	2343	1037	879	635
puffinn	-	2505	2695	1390	1037	593
flatnav	2994	2451	1691	1724	892	1141
falconnpp	-	3844	2025	1555	1203	676
nsg(faiss)	1968	2928	3121	2114	1442	1111
vamana(svs)	2333	3776	3328	1828	879	1040
hnswlib	3208	3174	3689	2339	1097	905
mrpt	2413	3758	4568	3252	1984	1973
symphonyqg	2844	4211	5812	2226	2055	1637
ngt-qg	3708	7405	7715	5829	3586	2341
ngt-onng	9161	11131	11365	6968	10474	10052

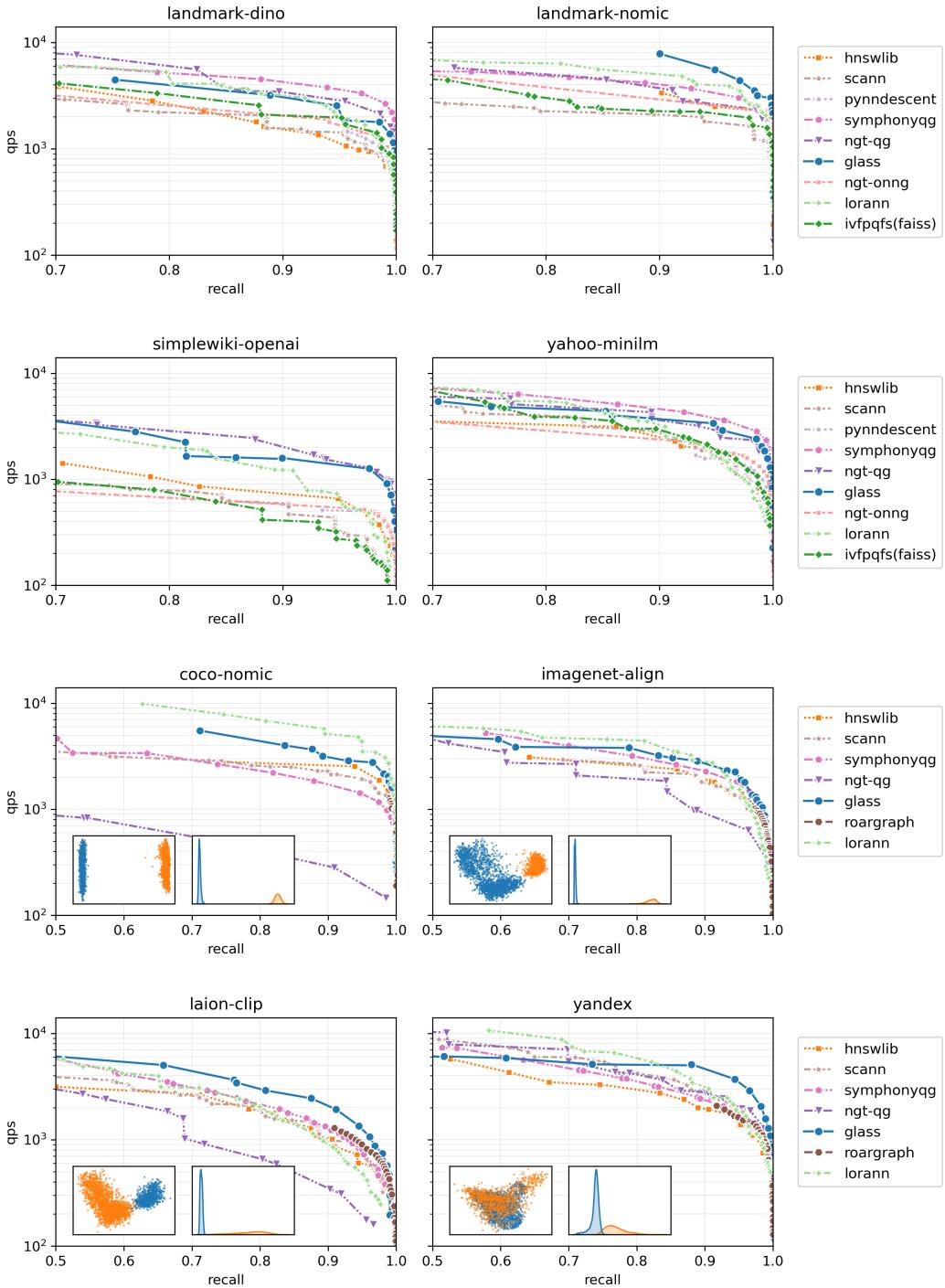
Comparing index construction times, it is worth noting that the hyperparameters in VIBE have been tuned for optimal query throughput, and index construction times may vary significantly depending on hyperparameter choices. We also note that, in general, the index construction of clustering-based and tree-based methods parallelizes more easily than the index construction of graph-based methods.

F Additional results

F.1 Selected algorithms

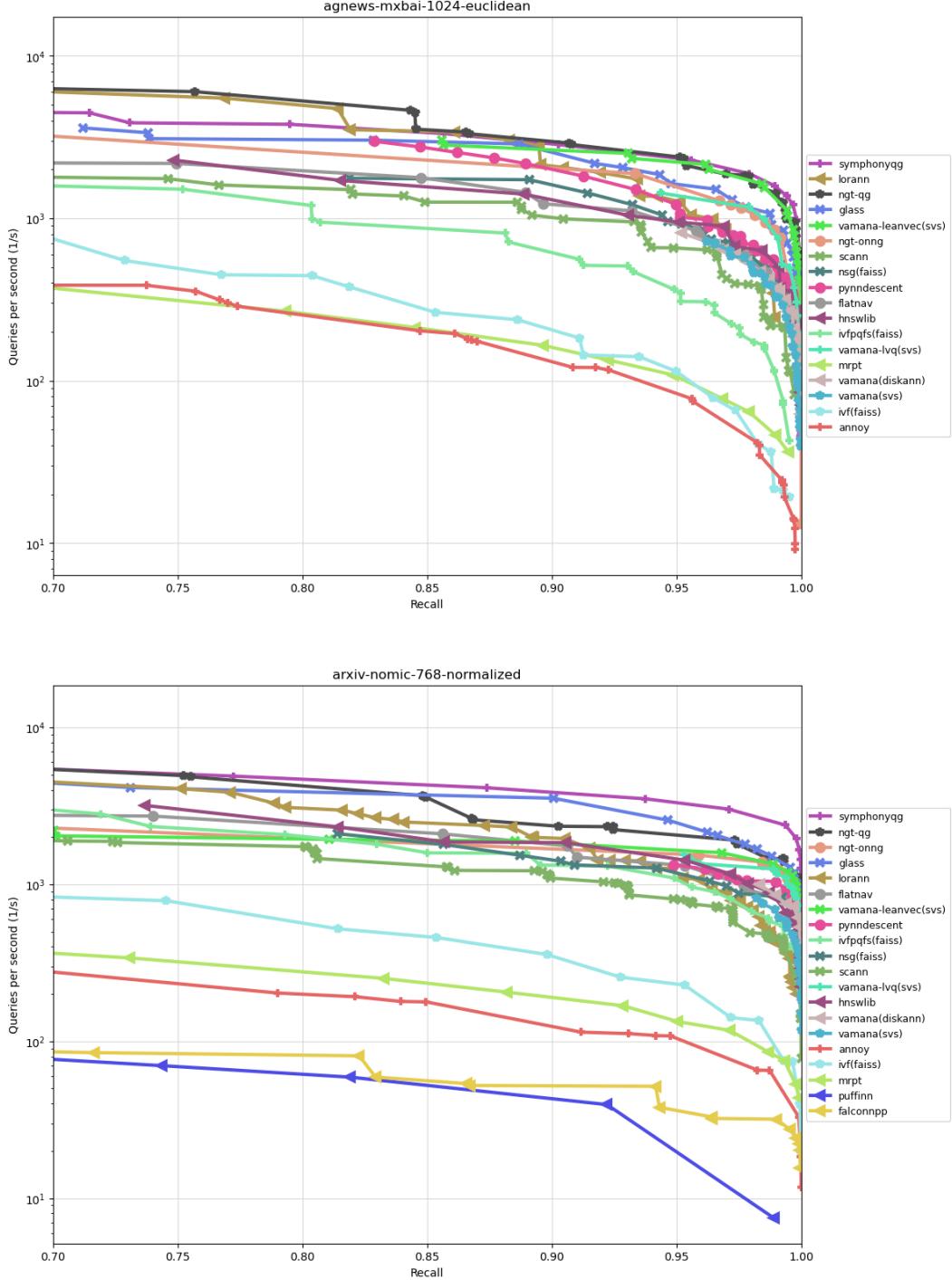
In this section, we present benchmarking results for all datasets using selected algorithms. For discussion, see Sections 5.2 and 5.3.





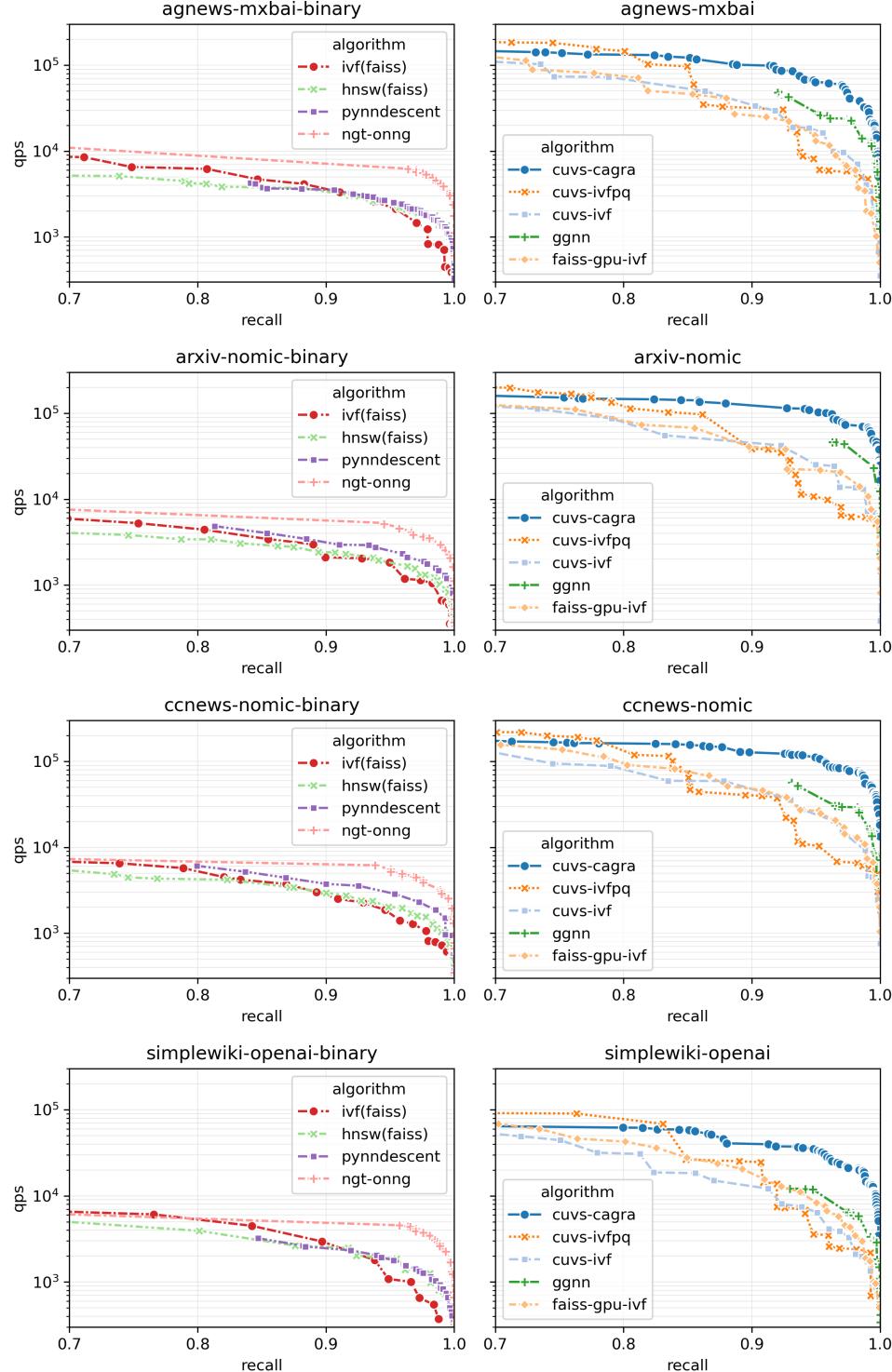
F.2 All algorithms

In this section, we present benchmarking results for all algorithms on two datasets. For full results on all datasets, we refer to our website: <https://vector-index-bench.github.io/>



F.3 Binary embeddings and GPU algorithms

In this section, we present additional results on achieving a higher throughput through binary embeddings and GPU algorithms. For details on the GPU algorithms, see Appendix D. For discussion, we refer to Section 5.2.



G OOD performance gap

In this section, we present the gap between the performance of in-distribution and out-of-distribution queries on the rest of the text-to-image datasets. For discussion, we refer to Section 5.3.

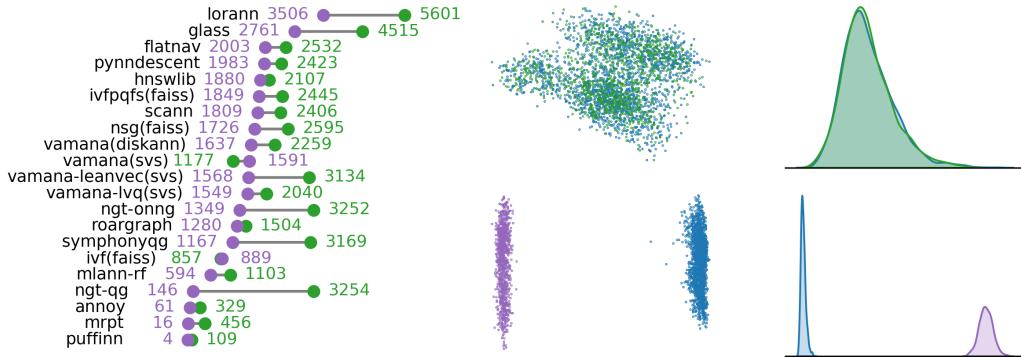


Figure 12: OOD performance gap on coco-nomic-768-normalized

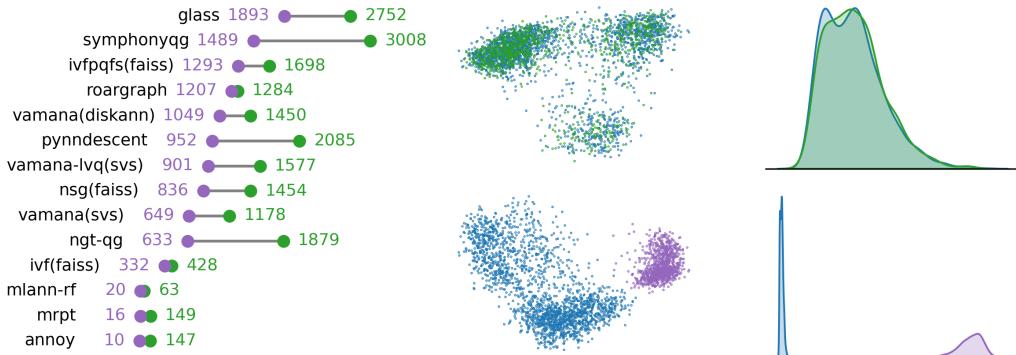


Figure 13: OOD performance gap on imagenet-align-640-normalized

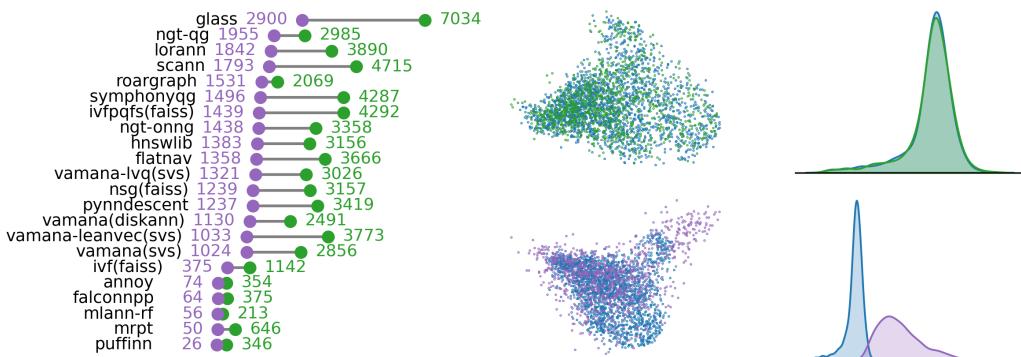


Figure 14: OOD performance gap on yandex-200-cosine